

Metabolomic Profiling of Nutrient Solutions for Characterization of Human Gut Microbiota

by

Cheng-Hui Sandi Yen

A thesis

presented to the University of Waterloo

in fulfillment of the

thesis requirement for the degree of

Master of Science

in

Chemistry

Waterloo, Ontario, Canada, 2014

©Cheng-Hui Sandi Yen 2014

AUTHOR'S DECLARATION

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

I understand that my thesis may be made electronically available to the public.

Overview

Studying the metabolite output, the “exometabolome,” of bacterial communities provides insight into the metabolic interactions occurring within that ecosystem. As such, a metabolomic approach can be used to gain further understanding of the human gut microbiota. As the microbes that reside in the human gastrointestinal (GI) system, the ecosystem has extensive impacts on the human body. *In vitro* systems for culturing gut microbes circumvents the ethical and technical challenges of *in vivo* studies, while, at the same time, lends itself to metabolomic investigations. The nutrient solutions in which microbial communities representative of the human gut microbiota are cultured contain a wealth of metabolomic data. In this thesis, metabolomic characterization is used to validate an *in vitro* gut-mimicking bioreactor system. This provided support for the ability to establish standardized protocols on preparing, culturing and reporting on *in vitro* gut microbiota. Furthermore, the metabolomic approach to studying gut microbes was used in applications relevant to clinical and industrial research. The findings from this thesis demonstrate the application of a metabolomic approach to studying the human gut microbiota in a single-stage continuous stirred tank reactor (CSTR) system. More importantly, the research here exemplifies the ways in which metabolomic insight can be translated into deciphering the complex interactions of the human gut microbiota.

Acknowledgements

This thesis would not be possible without the help and support of my mentors and peers. My thanks to Professor Marc Aucoin, who challenged me to ask the tough and relevant questions and in doing so, fostered a passion for science (and even engineering) within me. And to Dr. Emma Allen-Vercoe, who brought forth a wealth of knowledge and enthusiasm, in equal parts. I am grateful for both Marc's and Emma's mentorship throughout my degree. Thank you to my committee members, Dr. Guy Guillemette and Dr. Thorsten Diekmann for contributing time and energy into making this thesis the best that it could be.

I also extend my acknowledgements of my colleagues at the Aucoin lab. My deepest thanks go to Stan Sokolenko; thank you for your patience and your generosity. The skills you have taught me directly formed my work presented here and they will only continue to expand. Beyond this, it is impossible for me to overestimate the value of your friendship to me, as you have provided endless support to me, both in and outside of the lab. Jann Ang, Megan Logan, Steve George, Altamash Jauhar, Bhavik Manocha, Eric Blondeel, Valya Malenkov, Jian Xiong, and Biljana Todorovic, you have become my second family over the past few years. I am thankful to each of you for creating such a memorable and enjoyable workplace.

Further appreciation is extended to members of the Allen-Vercoe lab, without whom I would not have this thesis. Thank you to Julie McDonald, Erin Bolte, Kathleen Schroeter, Ian Brown and Chris Ambrose, for your time and commitment to the partnership between our two labs. As well, to Dave Chang of Chenomx Inc., who generously guided me in my introduction to the field of metabolomics. Lastly, thank you to my family who have been with me every step of the way. And to Nick, who is a constant to me.

Dedication

This thesis is dedicated to my parents, Mom and Dad, and to my siblings, Danny, Judy and Anita.

Table of Contents

AUTHOR'S DECLARATION	ii
Overview	iii
Acknowledgements	iv
Dedication	v
Table of Contents	vi
List of Figures	xiv
List of Tables.....	xvi
List of Abbreviations.....	xviii
Chapter 1 Introduction	1
Chapter 2 Literature Review	5
2.1 The Metabolomic Approach.....	5
2.1.1 Mass Spectrometry	6
2.1.2 Nuclear Magnetic Resonance Spectroscopy.....	6
2.1.3 Principal Component Analysis	10
2.2 ¹ H NMR-Based Metabolomics.....	11
2.2.1 Metabolomic study of Cell Culture Processes.....	12
2.3 The Human Gut Microbiota	13
2.4 Metabolomic Study of the Human Gut Microbiota.....	16
2.5 Studying the Gut Microbiota using Model Organisms.....	17
2.6 Studying the Gut Microbiota <i>In Vitro</i>	19
2.6.1 Batch Fermenters.....	19
2.6.2 Single-Stage Continuous Culture Systems	20
2.6.3 Multi-Stage Continuous Culture Systems	22
2.6.4 Mimicking Physiological Characteristics.....	24

2.6.5	Challenges of <i>In Vitro</i> Models.....	24
Chapter 3	Materials and Methods	26
3.1	NMR Sample Preparation.....	26
3.2	NMR Spectra Acquisition and Processing	26
3.3	Metabolite Profiles	30
3.4	Statistical Analysis	30
3.5	Bioreactor Operation and Inoculation	30
3.6	Collection and Preparation of Fecal Inocula	33
3.7	Preparation of Defined Communities.....	35
Chapter 4	Metabolite Profiles of Media.....	38
4.1	Chapter Objective.....	38
4.2	Constructing Metabolite Profiles.....	39
4.2.1	Section Objective.....	39
4.2.2	Processing NMR Spectra.....	39
4.2.3	Building Metabolite Profiles	41
4.3	Treating Cell Culture Media with UV Irradiation against Adventitious Agents: Minimal Impact on CHO Performance	43
4.3.1	Journal Article Authorship	43
4.3.2	Justification of Original Work.....	43
4.3.3	Abstract	44
4.3.4	Introduction	45
4.3.5	Materials and Methods	47
4.3.5.1	Experimental design	47
4.3.5.2	Media irradiation	48
4.3.5.3	Verification of UV fluence.....	48

4.3.5.4	Cell culture	51
4.3.5.5	Cell count.....	51
4.3.5.6	NMR spectroscopy and metabolite profiling.....	52
4.3.5.7	Antibody titer.....	52
4.3.6	Results	53
4.3.6.1	NMR metabolite analysis	53
4.3.6.2	Cell culture	56
4.3.6.3	Protein titre	56
4.3.7	Discussion.....	58
4.3.8	Conclusion.....	59
4.3.9	Acknowledgements	60
Chapter 5	Standardization of Fecal Water Samples and Analysis	61
5.1	Chapter Objective	61
5.1.1	Materials and Methods	62
5.1.1.1	Preparation of Fecal Water Samples	62
5.1.1.2	Statistical Analysis	64
5.1.2	Results	64
5.1.2.1	Fecal Water Analysis and Sample Preparation.....	64
5.1.2.2	Normalization of Metabolite Concentrations	66
5.1.2.3	Impact of NMR Spectra on Metabolite Profiles	69
5.1.3	Discussion.....	71
5.1.3.1	Trends in Metabolite Profiles	71
5.1.3.2	Normalization of Metabolite Concentrations	71
5.1.3.3	Impact of Spectral Limitations	72
5.1.3.4	Minimum Level of Feces Required in Fecal Water Samples	73

Chapter 6	Metabolomic Characterization of Gut-Mimicking CSTR	74
6.1	Chapter Objective	74
6.2	Gut-mimicking CSTR Effluent Sample Preparation and Analysis	75
6.2.1	Experiment Objective	75
6.2.2	Materials and Methods	75
6.2.2.1	Experiment Design	75
6.2.2.2	Sample Collection and Handling	77
6.2.2.3	Statistical Analysis	77
6.2.3	Results	78
6.2.3.1	Profiling Variability	78
6.2.3.2	Processing Methods	79
6.2.3.3	Sources of Effects Impacting Metabolite Concentration	79
6.2.3.4	Effect of Ultracentrifuge Speed	80
6.2.3.5	Comparison of Replicate Experiments	83
6.2.4	Discussion	86
6.2.4.1	Profiling Variability	86
6.2.4.2	Pre-Processing Treatments	86
6.2.4.3	Effect of Ultracentrifuge Speed	86
6.2.4.4	Comparison of Replicate Experiments	88
6.3	Effects of Varying Inoculum Biomass Proportions	89
6.3.1	Experiment Objectives	89
6.3.2	Materials and Methods	89
6.3.2.1	Experiment Outline	89
6.3.2.2	Sampling and Storage	91
6.3.2.3	Statistical Analysis	91

6.3.3	Results	92
6.3.3.1	Profiling Variability.....	92
6.3.3.2	Replicate Runs.....	96
6.3.3.3	Metabolite Profiles	99
6.3.3.4	Defined Cultures vs Their Reciprocals.....	102
6.3.4	Discussion.....	108
6.3.4.1	Profiling Variability.....	108
6.3.4.2	Replicate Runs.....	108
6.3.4.3	Metabolite Profiles	109
6.3.4.4	Defined Cultures vs Their Reciprocals.....	110
Chapter 7 Metabolomic Analysis of Human Fecal Microbiota Propagated <i>In Vitro</i> in a CSTR		
System 111		
7.1	Chapter Objective.....	111
7.2	Metabolomic Analysis of the Human Gut Microbiota: Comparison of Feces-Derived	
	Communities and Defined Mixed Communities.....	112
7.2.1	Journal Article Authorship	112
7.2.2	Justification of Original Work.....	112
7.2.3	Abstract	113
7.2.4	Introduction	114
7.2.5	Materials and Methods	116
7.2.5.1	Experiment overview.....	116
7.2.5.2	Bioreactor Operation and Inoculation	119
7.2.5.3	Preparation of Fecal Inocula and Defined Mixed Culture.....	119
7.2.5.4	Sample Collection and Processing	120
7.2.5.5	NMR Spectra Acquisition and Processing	120

7.2.5.6	Metabolite Profiles	121
7.2.5.7	Statistical Analysis	122
7.2.6	Results	123
7.2.6.1	Multivariate Analysis	123
7.2.6.2	Metabolite Profiles	129
7.2.7	Discussion.....	134
7.2.7.1	Feces-Derived vs Defined Community Cultures.....	134
7.2.7.2	Metabolite Profiles of Bacterial Communities	135
7.2.7.3	Donor A.....	136
7.2.7.4	Donor B	137
7.2.7.5	Donor C	138
7.2.7.6	MET-2	138
7.2.8	Conclusion.....	140
7.2.9	Acknowledgements	141
7.3	Analysis of Perturbations to Human Fecal Microbiota Propagated <i>In Vitro</i>	142
7.3.1	Experiment Objective.....	142
7.3.2	Justification of Treatment Condition.....	142
7.3.2.1	Clindamycin	142
7.3.2.2	Vancomycin.....	143
7.3.2.3	Norepinephrine	143
7.3.2.4	Defined Bacterial Community.....	144
7.3.3	Materials and Methods	144
7.3.3.1	Overview of Bioreactor Setup.....	144
7.3.3.2	Bioreactor Perturbations.....	146
7.3.3.3	Sample collection and processing.....	147

7.3.3.4	Statistical analysis	147
7.3.4	Results	148
7.3.4.1	Antibiotic treatment.....	148
7.3.4.2	Defined community	149
7.3.4.3	Norepinephrine treatment	151
7.3.5	Discussion.....	154
7.3.5.1	Clindamycin treatment	154
7.3.5.2	Vancomycin treatment.....	155
7.3.5.3	Challenge of Fecal Culture with Defined Community	156
7.3.5.4	Challenge of Fecal Culture in Dysbiosis with Defined Community	156
7.3.5.5	Norepinephrine treatment	159
Chapter 8	Conclusions	161
8.1	Metabolomic analysis of Nutrient Solutions	161
8.2	Standardization of Reporting on Fecal Water Metabolites.....	161
8.3	Metabolomic Evaluation of <i>In Vitro</i> Gut Bacterial Communities.....	163
8.3.1	Sample Preparation Procedure for Metabolomic Analysis.....	163
8.3.2	Effects of Varying Inoculant Biomass Proportions	164
8.4	Metabolomic Characterization of <i>In Vitro</i> Gut Microbial Communities	165
Chapter 9	Recommendations	167
References	170
Appendix A	Supplementary Figures and Tables.....	190
A.1	Materials and Methods	190
A.2	Standardization of Fecal Water Samples and Analysis	197
A.3	Sample Preparation Procedure for Metabolomic Analysis.....	200
A.4	Effects of Varying Inoculant Biomass Proportions	205

A.5 Metabolomic Analysis of the Human Gut Microbiota: Comparison of Feces-Derived Communities and Defined Mixed Communities	207
A.6 Analysis of Perturbations to Human Fecal Microbiota Propagated In Vitro.....	209
Appendix B Calculations.....	211
Appendix C Programming.....	213

List of Figures

Figure 2-1. Anatomy of the human GI tract	15
Figure 3-1. Schematic of 1D-NOESY pulse sequence.....	27
Figure 3-2. NMR spectrum of typical gut-mimicking single-stage CSTR sample after phase, baseline and shim correction steps.	29
Figure 4-1. Effects of each spectral processing step on the NMR spectrum	40
Figure 4-2. Identifying and quantifying compounds in NMR spectrum of defined CHO media.....	42
Figure 4-3. Received UV fluence in commercial media measured by bioassay (MS2 bacteriophage) vs. target value.....	50
Figure 4-4. Compound concentration as a function of UV fluence.....	54
Figure 4-5. Comparison of cell culture performance on control and treated media	57
Figure 5-1. Metabolite profile overview.....	65
Figure 5-2. Normalized compound concentrations	68
Figure 5-3. Profiling confidence impacted for certain compounds	70
Figure 6-1. Compounds that demonstrate significant linear relationship with centrifugal force.	81
Figure 6-2. Effect of molecular weight on percent concentration decrease	82
Figure 6-3. Compounds that had significant experiment bias based on MLR model.....	84
Figure 6-4. Experiment overview	90
Figure 6-5. PCA score plot of triplicate samples.....	95
Figure 6-6. PCA score plot of metabolite profiles of all bioreactor samples	97
Figure 6-7. The breakdown of compound classes in metabolite profiles of each defined culture.....	100
Figure 6-8. Compounds that reach the same level by Day 20 across a MET-2A culture and its reciprocal (MET-2B)	103
Figure 6-9. Common point reached in only one run for MET-2A and MET-2B.	104

Figure 6-10. Compounds that reach the same level by Day 20 across a MET-3A culture and its reciprocal (MET-3B)	105
Figure 6-11. Common point reached in only one run for MET-3A and MET-3B	106
Figure 7-1. Experiment overview	117
Figure 7-2. Principal component analysis score of metabolite concentration profiles	125
Figure 7-3. Metabolite profile comparison different donors.	126
Figure 7-4. Metabolites as identifying markers for metabolite profiles.	133
Figure 7-5. Experiment overview	145
Figure 7-6. Treatment effects on metabolite profiles	152
Figure A1. NMR Spectra of fecal water at various concentrations	197
Figure A2. Profiling confidence impacted for certain compounds.....	198
Figure A3. Significance of terms in each iteration of model simplification.....	200
Figure A4. Evaluation of profiling precision based on relative standard deviation of compound concentrations observed in single-step filtered samples from Experiment 2.	201
Figure A5. Density plot of the relative standard deviation of triplicate samples.	205
Figure A6. Profile overview of MET-2A and MET-3A and their reciprocal cultures, MET-2B and MET-3B.....	206
Figure A7. Kernel density plot of mean concentrations of untreated samples	207
Figure A8. PCA plots of metabolite profiles of donor A, B, C and MET-2 cultures	208
Figure A9. PCA model of metabolite profiles of Donor A, B, and MET-2, as well as profiles of each respective community culture after various treatments.....	210

List of Tables

Table 3-1. Composition of growth medium.....	32
Table 3-2. Fecal sample donors.	34
Table 3-3. Defined Communities and Their Parent Communities.....	37
Table 4-1. Changes in compound concentration per 100 mJ/cm ² of UV dose.....	55
Table 5-1. Fecal concentration of fecal water samples.....	63
Table 6-1. Sample treatment conditions.	76
Table 6-2. Absolute percent change due to experiment bias	85
Table 6-3. Compounds with RSD in replicate samples > 15%.....	94
Table 6-4. Compounds that differed significantly between duplicate runs.	98
Table 6-5. Proportion of Compound Classes, weighted to concentrations.....	101
Table 6-6. Compounds that do not reach a common point between a culture and its reciprocal by Day 20 in either of the duplicate runs.	107
Table 7-1. A list of cultured bacterial isolates from the feces of donor A included in the MET-2 defined community.....	118
Table 7-2. Compounds that differentiated donor-derived cultures and defined cultures as tested by the Mann-Whitney test at the 95% confidence level.....	127
Table 7-3. Metabolite profile diversity of each bacterial culture. Proportion of each compound class (%) weighted to the compounds' concentration.....	130
Table A1. Cultured bacterial isolates for MET-1.....	190
Table A2. Cultured bacterial isolates for MET-2.	192
Table A3. Cultured bacterial isolates for MET-2A, MET-2B, MET-3A, and MET-3B.	194
Table A4. Compounds that did not demonstrate significant linear trend with changes in centrifugal force, based on simplest MLR model.....	202

Table A5. Compounds that demonstrated the same trends in concentration change with ultracentrifuge speed in both Experiment 1 and Experiment 2.	203
-------------------------------------------------------------------------------------------------------------------------------------------------------------	-----

List of Abbreviations

¹³ C	Carbon
1D- ¹ H NMR	One dimensional proton nuclear magnetic resonance
1D-NOESY	One dimensional proton nuclear Overhauser effect spectroscopy
¹ H	Proton
BCFA	Branch chain fatty acid
BHI	Brain heart infusion
CHO	Chinese hamster ovary
Clm	Clindamycin
CMV	cytomegalovirus
CQA	critical quality attributes
CSTR	Continuously Stirred Anaerobic Bioreactor
DGGE	Denaturing gradient gel electrophoresis
DPI	Days Post-Inoculation
DSS	4,4-dimethyl-4-silapentane-1-sulfonic acid
FDA	Food and Drug Administration
FOS	Fructo-oligosaccharides
FTICR-MS	Fourier transform ion cyclotron resonance mass spectrometry
GC-MS	Gas chromatography mass spectrometry
GI	Gastrointestinal
HBM	Human Baby Microbiota
HMDB	Human metabolome database
HMP	Human Metagenomic Project
HTST	High-temperature/short time
LCFA	long chain fatty acid
LC-MS	Liquid chromatography mass spectrometry
M199	Media 199
MCC	Multi-stage continuous culture
MET	Microbial ecosystem therapeutic
MetaHIT	Metagenomics of Human Intestinal Tract
MS	Mass spectrometry
MVM	minute virus of mice
NE	Norepinephrine
NIST	National Institute of Standards and Technology
NMR	Nuclear Magnetic Resonance
PC	Principal component
PCA	Principal component analysis
Ppm	Parts per million
RED	Reduction equivalent dose

rpm	revolutions per minute
SCCS	Semi-continuous culture systems
SCFA	Short chain fatty acid
UV	ultraviolet
Van	Vancomycin
VFA	volatile fatty acids

Chapter 1 Introduction

The field of metabolomics focuses on the complete set of metabolites produced by an organism, referred to as the metabolome¹. As the summative product of an organism's biochemical processes, the metabolome produces a "snapshot" of that organism's metabolic state in that particular environment and at that particular time. A metabolomic approach could be taken to investigate the state of bacterial communities by observing the metabolite changes occurring in the nutrient solutions in which those microbes are grown. The metabolomic data of a bacterial community reveals information on both the metabolite production and nutrient consumption by the culture as a whole. As such, this type of characterization is directly applicable to understanding complex microbial consortia such as the human gut microbiota. As a consortium of bacteria residing in the large intestines of the human gastrointestinal (GI) tract, the metabolic interaction between the microbial community and the human host has direct implications on the physiology of both sides of the commensal relationship. On the one hand, there are a multitude of measurements that comment on the health and physiology of the human body, but assessment of the gut bacteria is not as straight forward. One means of studying the gut microbiota is by culturing the bacteria *in vitro* in a single-stage continuous stirred tank reactor (CSTR) system, and allowing the bacteria to grow as a community culture. There are many advantages in using such a system to study the gut microbiota, one of which is that it lends itself to metabolomic studies. In characterizing the metabolites in the spent media of the *in vitro* mixed culture, the human gut microbiota can be directly studied and subjected to experimentation.

The driving hypothesis of this thesis is that nuclear magnetic resonance (NMR)-based metabolomic data found in nutrient solution is an effective approach to studying the human gut microbiota *in vitro*. This overarching theme can be broken down into two subsequent sub-hypotheses. The first is that metabolomic data within spent media obtained from culturing a gut microbial

community can be used to evaluate the microbial changes imparted by *in vitro* growth, thereby validating the single-stage CSTR system used to mimic the human large intestines. This hypothesis will be tested by three main objectives. The first was to develop a standard method of reporting on fecal metabolomes. The second was to develop a protocol for preparing samples from the gut-mimicking bioreactor for NMR metabolomic analysis. The last objective was to determine the importance of the proportional bacterial makeup of fecal inoculum relative to the resulting *in vitro* community culture. The second hypothesis tested was that the metabolite profiles of gut bacteria communities grown in a single-stage CSTR system will sufficiently describe and compare the metabolic interactions taking place in the community culture, and define metabolic functions unique to each culture. Two main objectives were set to test this hypothesis. The first was to characterize feces-derived and defined communities that were representative of the human gut microbiota based on their metabolite signatures. The second was objective was to evaluate the communities' metabolic responses to different types of perturbations to the cultures' CSTR environment.

Establishing metabolite profiles for solutions that contained large numbers of metabolites is a difficult task with several challenges that can be overcome with training and practice. Therefore, the first goal was to demonstrate that NMR can be used to create metabolite profiles of a semi-well-defined solution such as serum-free cell culture media. In the first experimental chapter (Chapter 4), a targeted profiling approach is illustrated, followed by an experiment that demonstrates the applicability of the technique in investigating cell culture media. The experiment presented in this chapter examined the ability to create and use metabolite profiles of a chemically-defined media that underwent non-biological processing steps. More specifically, metabolite profiles were used to assess the chemical changes resulting from UV irradiation. As a nutrient solution, unspent defined cell culture media was an appropriate matrix to develop the analytical skills involved in building

metabolite profiles from complex ^1H NMR spectra. The examination shows an initial demonstration of targeted profiling techniques to build accurate and precise metabolite profiles.

Since fecal water has been the matrix of preference in metabolomic research on the human gut microbiota, Chapter 5 dressed the methodological concerns in reporting on fecal metabolomes. Even though fecal water is widely studied in the gut microbiota field²⁻⁴ and is used as the inoculant for gut-mimicking CSTR systems^{5,6}, there is a lack of standardization of the biofluid in terms of preparation for NMR-based metabolomic analysis. Therefore, samples of fecal water of various concentrations were examined by metabolomic analysis. This systematic evaluation of fecal water samples provided evidence to suggest a means of standardizing reports on fecal metabolome.

Next, the methodological approaches for studying human gut microbiota based on metabolomic descriptions of community cultures grown *in vitro* in an anaerobic CSTR was established (Chapter 6). Two experiments were conducted to characterize different aspects of propagating mixed cultures in a CSTR that mimicked the conditions of the human colon. First, attention was directed towards developing and validating an appropriate protocol for preparing samples collected from the gut-mimicking bioreactor. To this end, bioreactor samples were prepared for NMR-based metabolomic analysis using different methods. Comparing these different preparation protocols determined the effects of syringe filtration and ultracentrifugation on the metabolite profile of a given sample. The last experiment in this chapter studied CSTR-grown cultures that started from inoculants with various proportions of bacterial biomass. Given the limited standardization of the fecal inoculants used for *in vitro* community cultures, the impact of the inoculant on the *in vitro* culture was the focus of the investigation.

The final experimental chapter (Chapter 7) applied the metabolomic analyses developed in previous chapters to characterize the human gut microbiota. The first experiment presented compared

various bacterial communities, where these communities were either feces-derived, obtained from different individuals, or obtained through the culture of a defined community that was modelled on a fecal community. The objective of this study was to describe different gut bacterial communities based on their metabolic characteristics in a single-stage CSTR system. Having established a baseline metabolomic view of the gut microbiota, the second experiment in this chapter extended our analysis to characterizing the gut microbiota under various conditions. The objective of this investigation was to capture perturbations to the cultures' metabolome when the gut microbial community was subjected to different treatments. This study was a demonstration of how metabolomic investigation could be applied in a manner that would be relevant to clinical and industrial understanding of the human gut microbiota.

Experimental work pertaining to the preparation of fecal water, fecal inoculants and defined community inoculants; development and operation of the gut mimicking CSTR system; and collection and preparation of samples were all conducted in the Allen-Vercoe lab of the University of Guelph. The focus of this thesis was on the metabolomic implications of the experimental work and to conduct metabolic analyses on the collaborative projects.

Chapter 2 Literature Review

2.1 The Metabolomic Approach

Metabolomic research, originally rooted in plant⁷ and toxicology⁸ research, was developed as a complementary tool to functional genomics in order to incorporate molecular interactions within our understanding of biological systems^{9,10}. Since then, metabolomics has become an approach in its own right and has been applied in disease diagnosis¹¹, nutrition assessment¹² and microbial studies¹³. Metabolomic research is diverse in its application because it focuses on the complete set of low molecular weight metabolites produced by a given organism, referred to as the metabolome¹. Compared to other “-omes” such as the genome, transcriptome and proteome, the metabolome is most descriptive of the organism’s phenotype. The metabolome can contain a wide variety of molecules with different physical properties; from low molecular weight fatty acids and volatile organic acids to higher molecular weight amino acids, carbohydrates and lipids. Consequently, multiple analytical platforms are required to capture an entire metabolome. Furthermore, a comprehensive metabolome captures the metabolic state of the organisms at a given time, meaning that the metabolome studied is strongly influenced by temporal and spatial parameters of the experiment¹⁴. Therefore, a powerful metabolomic approach must begin with careful experimental design, followed by the use of analytical methods that are high-throughput and robust, and completed using appropriate statistical analyses of the collected dataset.

Given that no single analytical strategy can capture the metabolome as a whole, the choice in the platform selected involves trade-offs between advantages and disadvantages associated with each technique. There are two main analytical methods used to obtain metabolite data: mass spectrometry (MS) and nuclear magnetic resonance (NMR) spectroscopy.

2.1.1 Mass Spectrometry

Mass spectrometry (MS) is an analytical tool that identifies compounds by separating ions generated from organic or inorganic materials based on their mass-to-charge ratio^{10,15}. MS is usually coupled with a separation step such as gas chromatography (GC-MS) or liquid chromatography (LC-MS), though there are several other techniques that can be employed in conjunction with MS, such as Fourier transform ion cyclotron resonance (FTICR-MS)¹⁶. One of the main strengths of using MS is the sensitivity of the technique, where the lower limit of detection of compounds can be on the order of picamolar or nanomolar concentrations. MS has also been demonstrated to be robust, yielding reproducible identification and quantification of metabolites¹⁷. On the other hand, MS-based techniques require labour-intensive preparation steps that destroy the sample in the process of extracting, separating and detecting the metabolites¹⁸. Intrinsic to the separation process is a bias towards which metabolites are observed, therefore requiring additional derivatization steps to broaden the scope of compounds targeted¹⁹. For example, GC-MS targets volatile metabolites that are also thermostable up to 280 °C, but samples can be derivatized to convert compounds from non-volatile and thermolabile to volatile and thermostable before being analyzed by GC-MS^{20,21}.

2.1.2 Nuclear Magnetic Resonance Spectroscopy

NMR spectroscopy is another analytical tool by which to identify and quantify compounds. It functions by applying strong magnetic fields to induce atoms to spin at a precession rate proportional to the strength of the magnet applied. Within the externally applied magnetic field, about half of the nuclei will be in the lower energy “up” spin state and half in the higher energy “down” spin state. The gap between the two energy states corresponds to a specific frequency of electromagnetic radiation. A high-power radio frequency is applied to the sample to excite all the nuclei of a given type (i.e. ¹H). The sum of all the nuclei precessing at a specific resonant frequency form a net magnetization of the

sample. As the spins drop back from the high energy spin state to low energy spin state, the process of relaxation releases energy that can be observed by the NMR receiver as a signal, referred to as the free induction decay (FID) at that resonant frequency. Since an FID signal is a function of time, Fourier transform converts an FID signal to a function of frequency, which separates the signals into multiple peaks at their corresponding frequencies. The absorbed RF energy, which is exactly the energy gap between the two spin states, promotes nuclei in the low-energy spin-state to a high-energy spin-state. The energy emitted during the subsequent relaxation process from high-energy back to low-energy spin states is what is detected as that atom's signal. The spectrometer can use the characteristic frequency of a given nucleus, such as a proton (^1H) or carbon (^{13}C), and obtain information on all ^1H or ^{13}C atoms in a sample⁴.

One of the advantages of using NMR spectroscopy is that structural information of compounds found in the sample are also relayed through the resonant frequency. An atom's structural environment impacts its resonance frequency because small magnetic fields from surrounding atoms within a molecule influences the external magnetic field in the immediate surroundings of that nucleus. These variations to the external magnetic field, localized to individual nuclei, is described as the chemical shift, and is influenced by molecular symmetry and the types of electronegative atoms nearby, which have consistent and predictable effects on its surroundings. Therefore, the chemical shift of an atom, measured in parts per million (ppm) conveys structural information about its associated molecule and is reflected based on the location of the signal along the spectrum⁴. Additionally, nuclei in close proximity (within three bonds) affect the external magnetic field experienced by neighbouring nuclei, an effect called spin-spin splitting. Therefore, neighbouring nuclei mutually alter each other's resonance frequencies and alter the magnetic field experienced. Spin-spin splitting is represented in the spectrum by "splitting" the resonance peak corresponding to

the nucleus of interest into multiple peaks, set apart by J Hz, where J is the coupling constant. The pattern of peak splitting depends on the number of neighbouring nuclei and the intensity of the split peaks is determined by Pascal's triangle⁴.

NMR has also been demonstrated to be a robust analytical tool that yields reproducible results suitable for metabolomic analysis of biofluids^{22,23}. Furthermore, the sample preparation process for NMR spectroscopy is non-destructive and simple as the preparation procedure only involves adding an internal standard compound in deuterated solvent to the sample to control for drift in the magnetic field and for downstream quantification. Another benefit of using NMR spectroscopy for obtaining metabolomic data is the non-selective nature of the technology. Since NMR utilizes the quantum properties of molecules, the resonant frequency of all atoms are affected by the external magnetic field. The NMR spectrometer can target the characteristic frequency of a particular type of nucleus to obtain the resonance frequency information of all metabolites that contain the targeted atom within the sample. Most applications of NMR utilize ^1H or ^{13}C atoms to relay the metabolite make-up of a sample. This is because the sensitivity of NMR spectroscopy is directly related to the abundance of the targeted atom as well as the strength of that atom's nuclear magnet. For example, the reduced sensitivity of NMR to ^{13}C signals in comparison to ^1H is compounded by the fact that ^{13}C exists at much lower natural abundance (about 1.1%) and has a weaker nuclear magnet than ^1H . Since the sensitivity of the platform is dependent on the abundance of nuclei, NMR spectroscopy has relatively weak sensitivity. While there are hardware solutions to improve signal intensity and sensitivity, the lower limit of detection using NMR spectroscopy is still on the order of micromolar concentrations²⁴.

Just as there are many different throughput and separation techniques used to enhance various attributes of MS, there are many different NMR experiments that adjust the spectral acquisition

process to augment particular characteristics of NMR as a platform for obtaining metabolomic data. For example, there are experiments that aid with water suppression^{25,26}, improve signal de-convolution^{27,28}, or high-throughput interpretation of NMR spectra^{29,30}.

1D-¹H-nuclear Overhauser effect spectroscopy (1D-NOESY) pulse sequence has become the NMR experiment of choice when it comes to performing NMR-based metabolomic studies²³. Inherent to its popularity, one of the main benefits of using 1D-NOESY is that, as an established procedure, the results obtained from the pulse sequence can be consistently compared with other studies using the same experiment. Moreover, the widespread use of 1D-NOESY in metabolomic research has led to development of downstream processing tools, such as Chenomx NMR Suite (Chenomx Inc., AB), that facilitate the interpretation of the acquired spectra³¹. One factor that likely contributed to the ubiquitous use of 1D-NOESY for metabolomic analyses is the compatibility of most modern NMR spectrometers with the hardware requirements of the experiment. In terms of the spectroscopic advantages of the 1D-NOESY, it is an experiment that requires minimal optimization while still conferring adequate suppression of solvent signals³². 1D-NOESY exploits the difference in T_1 relaxation time of the solute and solvent in order to discriminate the two components of the matrix. Another aspect of 1D-NOESY that contributes to the successful suppression of the solvent signal is by using a series of pulse sequences to select against resonances from parts of the sample that did not experience the pulses at full efficacy or experienced non-homogenous fields. The third mechanism of 1D-NOESY that facilitates selection against unwanted signals is through the use of phase cycling. Phase cycling is a complex mechanism as its effects are embedded within many aspects of the experiment. Generally, signals of interest are “cycled through” various phases, so it is reversed relative to the phase of the NMR receiver. This way, external noise source contaminating a signal are subtracted while the signal of interest is added when alternately reversing the direction of the desired

signal and the NMR receiver. Given the quality of solvent suppression and the ease of integrating the resultant spectra into downstream analyses, this experiment can be appropriately applied to metabolomic studies.

2.1.3 Principal Component Analysis

The rich information contained within complex metabolomic data can be extracted using appropriate statistical tools. Multivariate analyses have been applied to interpret quantitative and qualitative annotations of metabolites, such as compound concentration and compound identity^{33,34}. Principal component analysis (PCA) is one such analysis that provides an overall view of the data structure, revealing trends and groupings within the data. This pattern recognition technique allows visualization of datasets that have many correlated variables associated with each observation. Within metabolomics, the metabolites are the variables associated with each observation, and an observation is a metabolomic spectra (NMR or otherwise). PCA is able to create a model based on the samples and their respective metabolite profiles by transforming the correlated variables into a reduced set of new, independent variables, called principal components (PCs)³⁵. Mathematically, the original data is represented as the covariance measured between every variable to form a covariance matrix. Once the variable relationships are captured in terms of their covariance, the eigenvectors of the covariance matrix represent score coordinates and the eigenvalues of the covariance matrix represent the loading values. The interpretation of these new PC variables in the context of metabolomics is that the score coordinate values illustrate how the spectra are related to one another, and the loading values depict how the metabolites contributed to the PCs. Graphically, when spectra are located along a given PC, it demonstrates that they are related along that axis, while the coordinate value of a metabolite in the same PC space illustrates how that metabolite relates to those spectra. Loading plots illustrate how the metabolites relate to those spectra³⁶. The ability to visually detect patterns of similarity and

differences within metabolite profiles is the main strength of PCA and has been the driving force for its popularity in metabolomic analysis.

2.2 ^1H NMR-Based Metabolomics

As mentioned earlier, the NMR spectrometer can use the characteristic frequency of different nuclei to relay information on the structure and concentration of metabolites in a sample. Since the sensitivity of the NMR is directly related to the abundance of that atom, monitoring the ^1H provides the greatest relative sensitivity as it is found at almost 100% natural abundance. For this reason, one dimensional proton NMR (1D- ^1H NMR) is often employed for obtaining metabolomic data in various fields, including cell culture research³⁷ and study of human biofluids³⁸. 1D- ^1H NMR spectra can provide quantitative identification of the spectral peaks through “targeted profiling,” a method whereby spectral peaks are individually matched with metabolite signatures obtained from a library of pure compound spectra³¹. This method of analyzing the spectra allows one to interpret the spectral signals, even when the peaks overlap, a characteristic often encountered in 1D- ^1H NMR spectra of complex solutions. Previous assessment of targeted profiling determined that the accuracy and precision of quantification is influenced by the combination of the degree of convolution affecting the metabolite’s signals and the magnitude of the compound’s concentration; compounds that are highly convoluted and found at low concentrations will have greater measurement error³⁹. Since measurement error in targeted profiling approaches are compound-dependent, peak-specific strategies have been developed to improve measurement accuracy^{39–41}. These strategies involve group-wide review of profiled spectra to establish confidence in each metabolite’s observation; generation of a master profiling list for consistent metabolite identification; or performing a validation experiment to define quantification accuracy within that sample type.

2.2.1 Metabolomic study of Cell Culture Processes

Within the framework of a single-cell organisms, the metabolome is composed of the endometabolome and the exometabolome, which refers the organism's intracellular metabolites and extracellular metabolites respectively¹³. Since the metabolome is regarded as the most accurate description of the organism's phenotype, it can be applied as an analytical tool in cell culture processes. For example, metabolomics has been used in culturing various cell lines, providing complementary data to genomic and proteomic measurements. In the field of industrial biotechnology, areas of focus for cell line metabolomics include phenotype classification, monitoring bioprocesses and ensuring consistent product quality. For example, metabolomic data have been used as a high-throughput screening tool for strain classification. *E. coli* strains containing mutations in the tryptophan metabolism pathway were identified based on the unique metabolite profiles associated with altered tryptophan pathways⁴². Discriminating mutant strains using metabolite profiles is directly relevant to industrial biotechnology since screening for and identifying cell line variants that strongly impact downstream products is highly desirable. Cell line metabolomics is also relevant to industrial processes as it can contribute to monitoring cell physiology during production. One example is in production of fuel ethanol using *Saccharomyces cerevisiae*. The environments imposed on the yeast cells during fermentation can illicit stress responses that could ultimately impact the end product⁴³. In such a case, the metabolic changes undergone by the cells in response to the harsh environments of fuel ethanol production were identified and contributed to determining the ideal growth conditions and media supplementation required for particular strains of yeast⁴⁴. Using metabolomic data for bioprocess monitoring also allows one to assess the production process. For example, NMR-based metabolomics was used to assess the scalability of producing a therapeutic protein in Chinese hamster ovary (CHO) cells⁴⁵. Aranibar et al found that biochemical signatures associated with differences in

production levels in a 5000-L and 7-L bioreactors were related to energy metabolism, like adenosine nucleoside phosphates and UDP-glucose, 1,6-glucose biphosphate and UDP-galactose. Another area of consideration in production of bio-therapeutics is assessing the nutrient feed provided to the cell cultures. Referred to as “fermentanomics,” assessing metabolite components in media^{46,47} could reveal information on product quality and cell culture productivity. Nutrient consumption and metabolite accumulation in the media can strongly influence cell culture performance^{47,48}. Therefore, rational medium designs paired with NMR-based monitoring have been executed to direct nutrient consumption in order to control product quality⁴⁶. One such example of this was the marked increase in the production of monoclonal antibodies in CHO cells by supplementing the media with additional methionine, tryptophan and tyrosine due depletion of those amino acids in conditions of excess glutamine and glucose.⁴⁷ The wide-ranging applications to which metabolomic data contributes to cell culture processes attests to the depth and versatility of the tool. Beyond analysing pure cell cultures, which provides metabolomic views of an entire organism, the metabolomic approach also has potential in deciphering whole biological systems, such as the human body.

2.3 The Human Gut Microbiota

The microbiota of the human gut is an intricate community that interacts symbiotically with the human host to influence health and physiology. It is estimated that the number of microorganisms residing within the human gastrointestinal (GI) system (Figure 2-1) is approximately 10^{14} cells, ten times that of the number of human cells in the body⁴⁹. Such a complex ecosystem of microorganisms, mostly bacteria⁵⁰, interact closely with its environment via genetic and metabolic signals. Over the past decade, there has been an increased interest in characterizing the human gut microbiota through analyses of such signals. On one hand, the establishment of the Human Metagenomic Project (HMP) Consortium and the Metagenomics of Human Intestinal Tract (MetaHIT) reflect some of the

collective efforts researchers have put forth to use genetic characterizations to identify the microbiota and the roles they play within the microbial ecosystem and with the human host^{51,52}. On the other hand, metabolomic studies have been used to explore and understand relationships involving the gut microbiota which are not encoded within the genome⁵³. Complementary to metagenomic studies, researching the metabolic interactions amongst the gut microbiota and with the host cells have contributed tremendously to establishing relationships between bacterial communities and with the human body, in healthy and disordered states⁵⁴.

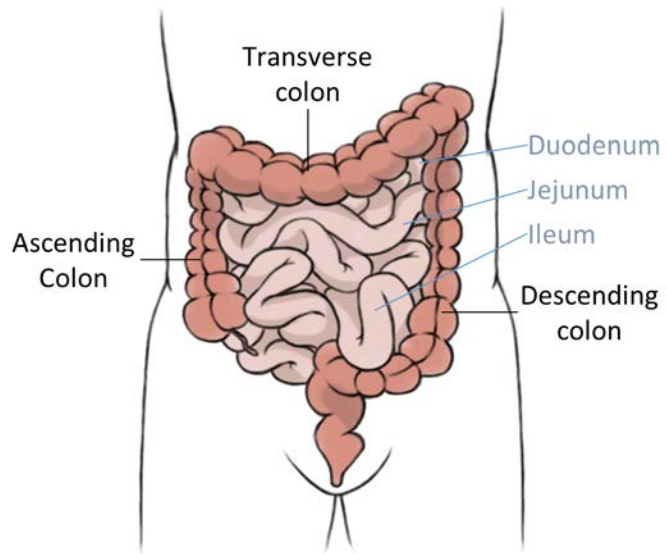


Figure 2-1. Anatomy of the human GI tract. Segments of large intestines labeled in black; segments of small intestines labelled in grey.

2.4 Metabolomic Study of the Human Gut Microbiota

When applied to bacterial communities, metabolomic approaches give insight into the metabolic activities engaged by the entire population. Within the context of the human gut microbiota, the mixed culture metabolome not only reflects the metabolic activities of the enteric bacteria, but it also represents the bacterial compounds to which the human body is exposed. Metabolomic studies of the human gut microbiota have been heavily focused on analyzing the compounds found in feces. Understanding gut bacteria through feces is an attractive approach as the procedure is non-invasive, yet still provides a wealth of information on the human-gut microbe relationship. Specifically, the study of the aqueous phase of feces, referred to as fecal water, has provided valuable information regarding host-bacterial changes under different conditions. The fecal metabolome was demonstrated to be an effective tool in assessing dietary effects, such as polyphenol-rich diets⁵⁵, daily apple intake⁵⁶, and red meat consumption⁵⁷. In clinical research, investigating the fecal metabolome has resulted in biomarkers for Crohn's disease⁵⁸ and colorectal cancer⁵⁹, and has been an effective diagnostic tool for ulcerative colitis³ and irritable bowel syndrome⁶⁰. The ability to obtain these findings via a non-invasive procedure is one of the main benefits of using feces to study the gut microbiota. Unfortunately, fecal samples have been shown to vary greatly between individuals and within an individual over time due to differences in gut physiology and variable diets, which in turn can affect properties of fecal samples, such as pH, moisture level and metabolite concentrations⁶¹. Therefore, one of the main drawbacks of examining stool samples is the heterogeneity of the matrix and the lack of standardization in preparing and analyzing the samples. There has been some effort in establishing optimised preparation procedures of fecal samples for the purposes of analyzing just the exometabolome⁶², or both the intra- and exometabolomes⁵⁵. However metabolomic data obtained from different preparation protocols and different batches of feces still cannot be performed without a

means of standardizing the sample's content. Interpretation of fecal metabolomes is further complicated by the fact that it is a reflection of host-microbe co-metabolism, where the fecal metabolome represents the cumulative effects on the host and the gut microbiota. Reproducibility of the metabolomic data is another obstacle in fecal metabolome research. The exact composition of the gut microbiota has been shown to vary with host age, nutrition, behaviour and health^{2,54,63}, so controlling for such experimental factors or recreating the exact conditions of previous studies can be quite difficult. Therefore, a system that allows for control over the experimental parameters subjected to the gut microbiota would circumvent the challenges of using fecal samples to study the gut microbiota.

2.5 Studying the Gut Microbiota using Model Organisms

The use of model organisms such as mice and rats allow for greater control over the subjects' diet and environment. The use of model organisms also widens the scope of experimental conditions that could be researched when it comes to understanding gut-host interactions, as some treatment conditions, such as antibiotic exposure, would have ethical limitations for human studies⁶⁴. For example, the administration of imipenem/cilastatin sodium, a broad spectrum β -lactam antibiotic to Wistar rats housed in a specific-pathogen-free environment was performed to assess its effects on the gut microbial-mammalian co-metabolism⁶³. Examination of the metabolome of both urine and feces demonstrated that profound changes to the host and microbe metabolism occurred rapidly (within the first day of treatment), while recovery from antibiotic treatment was gradual. The main metabolites that underwent large fold changes with treatment included aromatic benzene compounds, carbohydrates, short chain fatty acids (SCFA) and compounds of the tryptophan metabolism pathway. The metabolic perturbations observed were attributed to impacts on the gut microbiota rather than to host metabolic pathways since the antibiotic used had poor systemic absorption. At the same time,

host-microbial co-metabolism was still observed as there were changes to metabolites that are found in both mammalian and microbial physiology. Given that the gut microbiota are different across different species of model organisms and from humans², the applicability of studies performed *in vivo* can be enhanced through “humanizing” the model organisms’ gut microbiota. Germ-free or gnotobiotic mice, which are devoid of or have defined gut microbiota, can be inoculated with microbial consortia derived from human donors. Using this model, Martin et al. introduced human baby microbiota (HBM) to germ-free mice and examined the effects introducing *Lactobacillus paracasei* as a probiotic and a two dietary regiments (glucose-lactose mix and a galacto-oligosaccharide-rich prebiotic mix); all the while comparing them to conventional and conventionalized mice^{54,65}. Generally, they found that HBM mice had different levels of bacteria compared to the HBM mice that were supplemented with a probiotic, which also had different levels of bacteria than conventional and conventionalized mice. On the other hand HBM mice given a probiotic had similar levels of gut microbes as HBM mice supplemented with both a prebiotic diet and probiotic treatment, though the HBM mice subjected to the prebiotic diet had greater levels of bifidobacteria over time. In examining the fecal metabolic profiles, specific metabolic features could be attributed to the HBM community, the probiotic treatment and the two dietary regiments (basal glucose-lactose mix and prebiotic feed). Many of these identifying metabolic signatures involved changes to SCFA levels, such as acetate, butyrate and propionate, amino acid changes and also compounds associated with energy metabolism, such as succinate, glucose and lactose. Furthermore, there was a time-scale component to the changes observed, where the metabolite profiles of conventional and conventionalized mice were stable over time, while the HBM mice under different diets diverged from each other over time. These experiments conducted by Martin et al. exercised the power of NMR-based fecal metabolomics to examine the human gut microbiota under modulated environmental conditions. As exemplified in the studies described, studying the fecal metabolome in

an *in vivo* model allows one to assess the host-microbial co-metabolism and retain a moderate level of control over the experimental parameters, such as diet, environmental exposure and pharmaceutical treatment.

2.6 Studying the Gut Microbiota *In Vitro*

2.6.1 Batch Fermenters

The gut microbiota can also be studied through *in vitro* models. Where there are many model designs, the basis of these *in vitro* models involve the growth of fecal microbiota in a fermentation reactor operating anaerobically under physiological conditions. As the complexity of the bioreactor designs increase, so does the resemblance of the model to biological conditions. Batch reactors consist of culturing fecal inoculum in a single vessel supplied with nutrients, though without replenishment. Since batch reactors grow the bacteria in a closed system, the growth medium strongly influence the growth of the community culture. In one experiment⁶, the microbial communities reached similar maximum cell densities at the same time (18 h) in three kinds of brain heart infusion (BHI) media that varied in carbohydrate concentration. The different kinds of media had varying ability to maintain the initial fecal microbiota as measured by denaturing gradient gel electrophoresis (DGGE). For example, the low-carbohydrate BHI medium had the greatest ability to maintain the inoculum culture as indicated by similar DGGE profiles, while the high-carbohydrate BHI medium resulted in a DGGE profile with considerably fewer bands. In general, the *Bacteroidetes* phylum was more dominant in *in vitro* cultures, whereas the *Firmicutes* phylum were less dominant in *in vitro* cultures. Batch reactors are relatively cost effective, and easy to operate⁶⁶. For these reasons, batch reactors are suitable for experiments geared towards high-throughput screening for specific metabolic functions⁶, such as bioconversion of plant lignans by gut microbes⁶⁷. Batch fermenters accommodate growth of fecal

microbes only on short-term basis (24-48 h) because extended growth of the complex community can lead to unwanted selection of certain bacterial strains as nutrients are quickly consumed and the environment in the batch system changes over time. For this reason, continuous reactors that allow a more constant environment by replenishing nutrients and controlling pH are favored to mimic natural physiological conditions.

2.6.2 Single-Stage Continuous Culture Systems

Single-stage continuous stirred tank reactor (CSTR) models or single-stage semi-continuous culture systems (SCCS) allow for sustained nutrient feeding to the fecal microbes, and permit longer culture durations on the scale of several weeks. CSTR and SCCS models have been validated to maintain metabolic activity close to that of the initial fecal inoculum, as measured by hydrolytic and reductive enzyme activities, production of volatile fatty acids, hydrogen, methane and carbon dioxide⁶⁸. Shifts in fecal bacteria being cultured in single stage CSTRs and SCCSs include the enrichment of bacteria in the phylum *Bacteroidetes* and diminishment of bacteria in the phylum *Firmicutes*, similar to changes observed in communities grown in batch reactors. These shifts were independently observed in several validation studies, though the extent of the population shifts differed across studies. For example, Allison et al⁶⁹ observed slight decreases in overall count of anaerobes, with slight increases in *Bacteroidetes* and decreases in *Firmicutes* genera like *Clostridia*, and *Eubacteria*, though increases in other *Firmicutes* like *Lactobacilli* and *Streptococci* were also observed over a 14-day period. Campbell et al⁶⁸ saw a decrease in the total number of microbes in their SCCS-cultured fecal communities. The fecal inocula contained seven different genera, with 81% being *Bacteroides* and *Eubacterium*. 31% of anaerobic colonies could not be subcultured. In comparison, the 21-day *in vitro* culture only had four different genera characterized, with 89% being *Bacteroides*, *Eubacterium*, and *Clostridium* spp and only 13% of isolated anaerobes could not be subcultured. In the CSTR system

used by McDonald et al⁵, fecal microbial communities became enriched in phylum *Bacteroidetes* and diminished in phylum *Firmicutes*, such as *Clostridium* (cluster XIVa), and *Bacilli* after 4-5 weeks in the *in vitro* model.

Single stage CSTR systems and SCCS reactors also lend themselves to studying the effects of prebiotics and probiotics on gut microbial communities. Prebiotics are nondigestible food ingredients that selectively stimulate the growth or activity of bacteria in the GI tract, with the purpose of beneficially impacting the host⁷⁰. Probiotics refers to live microbial supplements used to improve the intestinal microbial balance, with the aim of beneficially impacting the host⁷¹. Duncan et al⁷² used a range of carbohydrates to study substrate impact on microbial composition using 16S rRNA probes. The investigators then continued on to use these substrates to screen for highly competitive strains that had potential to be used as new probiotic strains. The total number of anaerobes remained the same with changing substrates, except when fed with pectin, which causes a global decrease in the number of anaerobes. Furthermore, specific bacterial populations were strongly impacted by certain substrates. Dhalia inulin increased *Eubacterium cylindroides* by 100-fold and *Ruminococcus bromii* and *Ruminococcus flavefaciens* by 10-fold. The change in bacterial composition was also reflected in the production of SCFA. The proportional make-up of acetate-propionate-butyrate metabolites showed that acetate production was greatest with pectin (75%-16%-8%) and that growth on amylopectin produced the highest proportion of propionate and butyrate (44%-42%-14%). The impact of different substrates on the ability of probiotic candidates to compete with resident gut microbiota varied, with inulin and amylopectin facilitating the survival of *Rosuria* strains. Another application of single-stage CSTR has involved the study of fructo-oligosaccharides (FOS) as a prebiotic to prevent the degradation of isoflavone genistein by the gut microbiome and preserve its anti-carcinogenic properties⁷³. Steer et al. found that FOS increased the total count of anaerobic

bacteria, where *Bifidobacterium* spp. and *Lactobacillus* spp. were enriched and *Bacteroides* spp. and *Clostridium* spp. were decreased. Additionally, isoflavone genistein degradation was detected to a lesser extent with FOS-supplemented media than the unmodified recipe. These two examples illustrate the applicability of using single-stage continuous reactor systems to investigate the metabolic functions of the gut microbiota beyond the screening capabilities of batch reactors.

2.6.3 Multi-Stage Continuous Culture Systems

Characterizing the entire large intestine as a continuous culture system is still a simplification of the organ's operation. In truth, only the ascending colon receives steady nutrients from the small intestines and performs peristaltic mixing of the colon contents. The transverse and descending colon experience restricted nutrient supply and no mixing. As such, the different sections of the colon – ascending, transverse and descending – can be represented *in vitro* by using multi-stage continuous culture (MCC) systems. There have been many types of reactor designs that aim to recreate the spatial, temporal, and biochemical aspects of the large intestine. Bearne et al. designed a two-stage reactor system, where the two vessels differed in pH⁷⁴, while Macfarlane and Gibson developed a three-stage continuous culture system. In the three-stage system, the ascending, transverse and descending colon are mimicked by connecting three vessels differing in pH, operating volumes and dilution rates⁷⁵. Compartmentalizing the different segments of the large intestines revealed that the bacteriological content of the three reactor resembled that of the large intestines in terms of total anaerobic counts, but differed from *in vivo* examples in terms of species make up. Population densities of *Bifidobacteria* and *Bacteroides* were similar to the human large intestines. But the levels of *Clostridium perfringens* were greater in the first and second reactors than the ascending and transverse colon, while *Escherichia coli* was greater in the third reactor than in the descending colon. These species variation from one reactor to the next also meant that the vessels were also different

metabolically. Carbohydrate fermentation occurred mostly in the first “ascending colon” reactor, while an increase in branch chain fatty acids (BCFA) and phenolic compound due to the fermentation of amino acids that was observed in the second “transverse colon” reactor and third “descending colon” reactor.

A more comprehensive representation of the human gastrointestinal (GI) system is the Simulator of the Human Intestinal Microbial Ecosystem (SHIME) system developed by Molly et al⁷⁶. In the SHIME system, five consecutive reactors varying in working volume and residence time are used simulate the entire gastrointestinal tract. The first two reactors mimic the small intestine, where the duodenum and jejunum are represented in the first reactor and the ileum is represented in the second reactor. These reactors received successive inoculations of nutrition suspension representative of human western diets. The last three reactors represented the three sections of the large intestine with each vessel receiving a single fecal inoculation. Enumeration of the bacteria revealed that the microbial compositions in the SHIME system were comparable to *in vivo* communities. Metabolic assessments were made through measurement of volatile fatty acids (VFA), headspace gases (CO₂, CH₄, and H₂) and ammonium. Based on these parameters, there were significant differences between the SHIME culture and *in vivo* values from literature. Acetic acid, and propionic acid were generally found at higher molar ratios and butyric acid was consistently detected at lower molar ratios in the SHIME simulator compared to *in vivo* literature values. The headspace gas composition and ammonia levels were consistent across the various media types used, but differed from reference values. This divergence could be attributed to the use of O₂-free N₂ to maintain an anaerobic atmosphere. Constantly sparging the headspace with N₂ has been shown to impact bacterial metabolism⁷⁷. Recently, there have been modifications made to existing single stage and MCC systems to create anaerobic conditions using only the gasses produced by the fecal bacteria^{78,79}.

2.6.4 Mimicking Physiological Characteristics

The continuous culture systems described thus far allow metabolites and microbes to accumulate in the reactors until steady-state is reached all the while having both compounds and microorganisms removed together. The drawback of this system is that *in vivo*, metabolite concentrations are maintained relatively low as they are absorbed by the host. The accumulation of metabolites in *in vitro* systems can inhibit growth of certain microbes. One design that has attempted to address these concerns, has incorporated peristaltic mixing, water absorption and metabolite absorption^{80,81}. In this design, a dialysis membrane enclosed in a tubular system has been used to remove bacterial metabolites and water. This, combined with flexible silicone membranes undulated by water pressure to simulate peristaltic mixing of luminal contents, further allows the reactor to be inoculated at high-densities despite the viscosity of the solutions. In this system, microbial stabilization was reached in 24 hours and the *in vitro* community was sustained stably for over three weeks. Validation and characterization of *in vitro* microbial diversity in this system^{80,82} has led to applications such as investigating the effects of various carbohydrates, such as inulin⁸³ and lactulose⁸⁴, on microbial composition and metabolic activity and attributing carbohydrate fermentation functions to specific microbial populations by using isotope-labelled substrates and isotope-labelled RNA probing^{85,86}.

2.6.5 Challenges of *In Vitro* Models

Despite the advances in developing bioreactor designs to mimic the human GI tract, there are still challenges and limitations to using *in vitro* models. With respects to physiological factors, *in vitro* studies of gut microbiota lack host-microbe interactions. As such, the impact of the microbial activities on the host organism and the host's feedback reaction to the microbes cannot be measured. Furthermore, direct comparisons of community cultures from different fermentation experiments, even those using similar models, cannot be made. Part of the variability is attributed to inter-

individual differences in the bacterial make up within the gut. Since the gut microbial ecosystem is individually unique, *in vitro* communities are inherently different as they are derived from feces collected from different sources. Therefore, inter-individual variability must be taken into account in experimental designs, ensuring that treatment or disease effects are greater than biological variability. In terms of engineering challenges, variable retention times in the bioreactors is another confounding factor that impairs cross-experimental comparisons. Child et al⁸⁷ found that retention time strongly influences substrate availability, and bacterial colonization, which in turn shapes the entire *in vitro* community. Some bioreactor model validation studies take this factor into account and characterize its effects^{75,81}. Another aspect of consideration is defining when steady-state has been established. The amount of time required to reach steady-state is highly variable, even within a single bioreactor design. For example, single-stage continuous culture systems have been observed to reach steady-state as early as 5 days to as late as 4-5 weeks after inoculation^{5,68,69}. Since the achievement of steady-state is not independent of retention time, it is expected that the time required to reach a stable microbial environment is variable across *in vitro* models. Rather, attention should be focused on establishing a definition of steady state, whether it be based on microbial composition, metabolic production, or functional activity, and confirm its achievement within a given model. Lastly, *in vitro* models could easily lend themselves to metabolomic studies. However, the large majority of the *in vitro*-based investigations have placed its emphasis on microbial composition and microbial activity. Considerations of metabolic functions were usually probed only to the extent of VFAs or SCFAs or gaseous products such as H₂, CO₂, CH₄ and NH₄. Holistic examination of the metabolome of *in vitro*-propagated communities could be the next avenue of characterization in the movement to understand the human gut microbiota.

Chapter 3 Materials and Methods

3.1 NMR Sample Preparation

Samples analyzed with NMR spectroscopy were prepared to contain 10% internal standard solution, which was composed of 5 mM 4,4-dimethyl-4-sialpentane-1-sulfonic acid (DSS) and 0.2% w/v sodium azide in 99.9% D₂O (Chenomx Inc., Edmonton, Canada). This internal standard solution served as a chemical shift indicator (CSI) for software analysis of the spectra, while at the same time inhibiting bacterial contamination of the samples. The sample-DSS solutions were vortexed and transferred to 5-mm NMR tubes (NE-UL5-7, New Era Enterprises Inc., Vineland, NJ) for scanning. Samples were prepared on the day of scanning, and kept on ice until ~10 min before being scanned to allow the sample to return to room temperature.

3.2 NMR Spectra Acquisition and Processing

The 1D-¹H NMR spectra were acquired at the University of Waterloo Nuclear Magnetic Resonance Facility using a using the first increment of a NOESY pulse sequence (1D-NOESY) on a 600 MHz Bruker Avance NMR spectrometer. The NMR spectrometer was operated using XWIN-NMR 3.5 software and a Triple Resonance Probe (TXI 600). The 1D-NOESY pulse sequence compatible with the Chenomx software used to analyze the NMR spectra involved a 1 s presaturation that consisted of a 10 ms relaxation delay with a 990 ms presaturation period for water suppression at an induced field strength/ γB_1 of ~80 Hz⁴¹. This was followed by two 90° high power pulses and a 100 ms mixing time and another 90° high power pulse before a 4 s acquisition period. The total recycle delay in the experiment was 5 s. A total of 4 steady state scans and 32 scans were performed to obtain NMR spectra for a given sample (Figure 3-1).

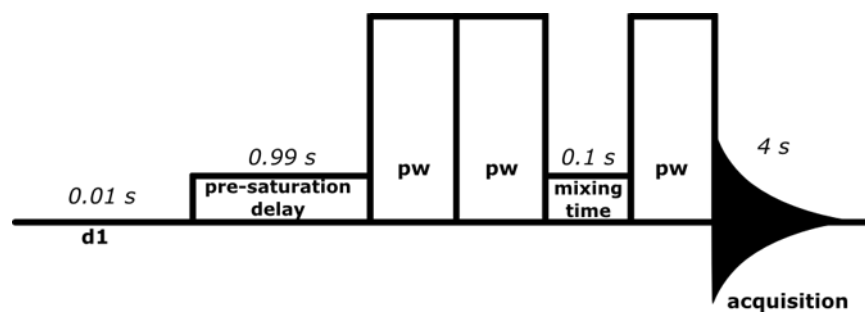


Figure 3-1. Schematic of 1D-NOESY pulse sequence. Figure credited to Chenomx application note on Investigating the criteria for accurate quantitative results with Chenomx NMR Suite⁸⁸.

Chenomx software was used to construct metabolite profiles from the 1D-¹H NMR spectrum (Figure 3-2). The corrected experimental spectra were assigned compounds by targeted profiling, superimposing compound signals in the built-in library to spectral peaks. The area of the spectrum signal relative to the reference compound was used for metabolite quantification. Profiled concentrations were divided by 0.9 to adjust for dilution with the internal standard solution during sample preparation for NMR spectroscopy. Sample spectra in each experiment were profiled according to scan order.

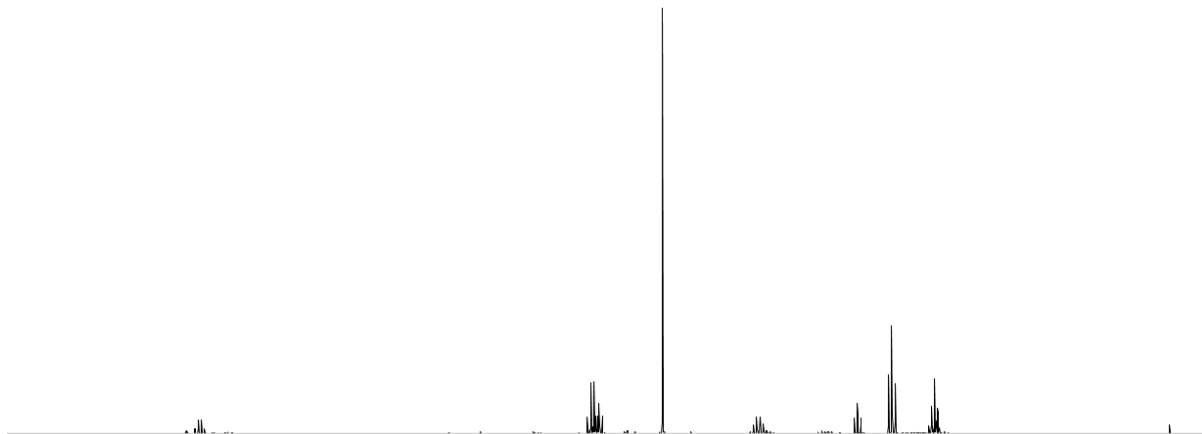


Figure 3-2. NMR spectrum of typical gut-mimicking single-stage CSTR sample after phase, baseline and shim correction steps.

3.3 Metabolite Profiles

The Chenomx software was also used to construct metabolite profiles from the 1D-¹H NMR spectrum. The corrected experimental spectra were assigned compounds by targeted profiling, superimposing compound signals in the built-in library to spectral peaks allowed for assignment of compound identity, and the area of the spectrum signal relative to the reference compound was used for metabolite quantification. Sample spectra in each experiment were profiled according to scan order. Once the metabolite profiles were assembled, each compound was assessed post hoc for confidence in the compounds' identification and quantification. Confidence in a compound was determined by the number of clusters in the compound's NMR signal, level of convolution with other compound signatures, and the signal-to-noise ratio of the compound. Only high-confidence compounds were included in subsequent analyses.

3.4 Statistical Analysis

All statistical analysis was performed using R, the statistical programming language⁸⁹. The following packages were used: RSQLite, RODBC, ggplot2, gridExtra, reshape2, plyr, and RColorBrewer and all base packages accompanying R.

3.5 Bioreactor Operation and Inoculation

Bioreactor operation and experimentation was performed by members of the Emma Allen-Vercos lab at the University of Guelph. Preparation and operation of the single stage, continuous stirred-tank reactor has been previously described^{90,91}. The Infors Multifors bioreactor system (Infors, Switzerland) was used throughout this body of work. 500-mL Multifors bioreactors (3 Multifors systems, Infors AG, Switzerland), with working volumes of 400 mL were operated as continuous

stirred tank reactors. Conditions were set to mimic the distal human gut (37°C, pH 6.9-7.0, gentle agitation, oxygen-free, and fed with a constant supply of mucin and insoluble starch substrates at a flow rate of ~400 mL/day). Maintenance of pH was automated with addition of 5% v/v HCl and 5% w/v NaOH. Vessels were maintained anaerobic by bubbling O₂-free N₂. The growth medium used in the bioreactor system was developed by the Emma Allen-Vercoe lab at the University of Guelph, and has been previously described^{90,91}, and is outlined in Table 3-1. Prior to inoculation, a sample of growth medium was collected from the bioreactor vessel to assess for bacterial contamination. The media sample was plated on fastidious anaerobe agar (Acumedia; Lansing, Michigan) supplemented with 5% defibrinated sheep's blood (Hemostat Laboratories; Dixon, California) and incubated in aerobic or anaerobic conditions, at 37 °C. The CSTR system used to culture the mixed bacterial communities was inoculated with either 25% (v/v) of 10% (w/v) fresh fecal slurry, or 12.5% (v/v) of defined culture suspension relative to the bioreactor working volume. Once inoculated, the culture was gently agitated. Medium feed was started 24 h post-inoculation.

Table 3-1. Composition of growth medium.

Reagent¹	Concentration (g/L)
Peptone water ^b	2
Yeast extract ^c	2
NaHCO ₃ ^b	2
CaCl ₂ ^a	0.01
Pectin (from citrus) ^a	2
Xylan (from beechwood) ^a	2
Arabinogalactan ^a	2
Starch (from wheat, unmodified) ^a	5
Casein ^d	3
Inulin (from Dahlia tubers) ^d	1
NaCl ^a	0.1
K ₂ H ₂ PO ₄ ^a	0.04
KH ₂ PO ₄ ^b	0.04
MgSO ₄ ^e	0.01
Bile salts ^a	0.5
Hemin ^e	0.005
Menadione ^a	0.001
L-cysteine HCl ^b	0.5
Porcine gastric mucin (type II) ^a	4

¹Chemical suppliers are indicated in superscripts: (a) Sigma-Aldrich (Oakville, Ontario); (b) Thermo Fisher Scientific (Ottawa, Ontario); (c) BD (Franklin Lakes, New Jersey); (d) Alfa Aesar (Ward Hill, Massachusetts); (e) BDH (Radnor, Pennsylvania). Chemical suppliers Adapted from McDonald et al, 2014⁵.

3.6 Collection and Preparation of Fecal Inocula

Collection of fecal samples and the preparation of fecal inocula was performed by members of the Allen-Vercor group as per the approval of the Research Ethics Board of the University of Guelph (REB#09AP011). The procedure was previously described and is summarized here^{90,91}. Fresh fecal samples were obtained from four healthy individuals (Table 3-2), none of whom had a recent history of antibiotic treatment prior to time of donation. Donors provided fresh fecal samples in a sealed sterile plastic container which was placed in an anaerobic chamber (atmosphere 90% N₂, 5% CO₂ and 5% H₂) within 5-10 min of defecation. Each fecal samples were processed into a 10% (w/v) fecal slurry by mixing 5 g of feces in 50 mL of growth medium for 1 minute with a stomacher (Tekmar Stomacher Lab Blender, Seward; Worthing, West Sussex, UK). Once homogenized, the fecal slurry was centrifuged (10 min, 175 x g) to remove large particulates, such as residual undigested food⁹². The remaining supernatant was used as fecal inoculum for the gut-mimicking bioreactors used in the studies.

Table 3-2. Fecal sample donors.

Donor	Sex	Age (yr)	Time elapsed between most recent antibiotic treatment and sample collection
A	Male	44	6 years
B	Female	26	9 months
C	Male	25	11 months
D	Female	42	8 years
E	Female	7.72	6 months

3.7 Preparation of Defined Communities

The defined community cultures were prepared and described by Schroeter, 2014⁹¹. These defined bacterial communities were referred to as Microbial Ecosystem Therapeutic (MET) cultures as they were originally developed as a community culture probiotic for therapeutic purposes. In short, pure bacterial isolates, identified by 16S rRNA gene sequencing, were obtained from fresh feces collected from two different donors: Donor D, a 41-year old female and Donor A, a healthy 44-year old male. 33 species were selected from the axenic fecal cultures isolated from Donor D to prepare a defined community, MET-1. Similarly, MET-2 was prepared from 33 species selected from Donor A. The breakdown of bacterial isolates found within MET-1 and MET-2 can be found in Table A1 and Table A2, respectively. Additionally, another 32 species from Donor D were selected to prepare the defined community, MET-2A. From the axenic fecal cultures isolated from Donor A, 32 species were selected to prepare the defined community, MET-3A. Once each of the bacterial strains were grown in pure culture using appropriate growth medium and conditions, an inoculum was created by combining biomass scraped from each culture plate and suspended in degassed 0.9% saline (37 °C). The final suspension volume for MET-1 was 30 mL while the final suspension volume for MET-2, MET-2A and MET-3A was 75 mL. The defined communities MET-2B and MET-3B were designed to be variants of MET-2A and MET-3A, respectively. MET-2B and MET-3B were prepared the same way as MET-2A and MET-3A, but differed from their respective parent cultures in that almost every isolate was added to the inocula in different proportions. The bacterial isolates for each of the four cultures and details of their biomass are detailed in Table A3. The same procedure for preparing defined community inoculants was repeated for the replicate experiment. Due to poor growth of four of the isolates in the pure cultures, the inocula for each of the defined communities had a different amount of bacteria for those poor-growth isolates as compared to the original experiment. These

differences are highlighted in Table A3. None of the defined community inocula included all four of the poor-growth isolates, so each defined community inocula only had 2-3 isolates that were different between the two replicate experiments. Given that the variant communities, MET-2B and MET-3B, had altered biomass in nearly every isolate, the deviations between the replicate experiments' inocula was considered to be negligible as any variations contributed by such differences would be indistinguishable from noise. A summary of the defined communities and the parent communities from which they were derived is shown in Table 3-3.

Table 3-3. Defined Communities and Their Parent Communities

Culture	Number of Bacterial species	Derived from Fecal Culture	Derived from Defined Culture
MET-1	33	Donor D	-
MET-2	33	Donor A	-
MET-2A	32	Donor D	-
MET-2B	32	-	MET-2A
MET-3A	32	Donor A	-
MET-3B	32	-	MET-3A

Chapter 4 Metabolite Profiles of Media

4.1 Chapter Objective

While the process constructing metabolite profiles is relatively straightforward, doing so in an accurate and precise manner requires technical refinement. In this chapter, the techniques involved in obtaining the metabolite profiles of nutrient solutions is demonstrated, first through a step-by-step illustration, then in an experiment that exemplifies the application of the technique for the purpose of metabolomic analysis. The experiment presented in this chapter examined the metabolites of a semi-well defined media for culturing Chinese hamster ovary (CHO) cells and investigated the effects of UV irradiation of that media on the metabolite profile of the media and on the media's ability to support the culture of CHO cells. The analysis of semi-defined, unspent cell culture media is an appropriate platform to introduce the NMR-based metabolomic techniques developed thus far to analysis of nutrient solutions.

4.2 Constructing Metabolite Profiles

4.2.1 Section Objective

The purpose of this section is to illustrate NMR spectra of solutions containing a large number of metabolites and demonstrate how metabolite profiles are constructed from these convoluted spectra using Chenomx NMR suite software.

4.2.2 Processing NMR Spectra

The raw NMR spectra were processed manually as described in Chapter 3. Figure 4-1A depicts the raw NMR spectrum (0.9-1.5 ppm) of a defined media for CHO cells, where the signal peaks are out of phase, and neither baseline correction nor shim correction have been applied. Once the spectrum has been phase corrected (Figure 4-1B), baseline corrections can be made by setting the baseline manually. Figure 4-1C exemplifies the same NMR spectrum but with phase and baseline corrections applied. Figure 4-1D and Figure 4-1E show how shim correction is applied throughout the spectrum using the internal shape indicator, DSS. The blue DSS template is overlaid with and adjusted to fit the spectral DSS signal at 0.0 ppm. Mismatches between the spectral DSS signal and the software DSS template is considered to be a spectral artifact produced by imperfections in the NMR instrumentation. Therefore, all the NMR spectrum signals are altered to match DSS (Figure 4-1E). Figure 4-1F shows the NMR spectrum (0.9-1.5 ppm) after all the NMR processing steps. Each of these steps aid in producing the ideal NMR spectrum, without altering signal peak area which will facilitate signal identification and compound quantification in subsequent steps of building a metabolite profile.

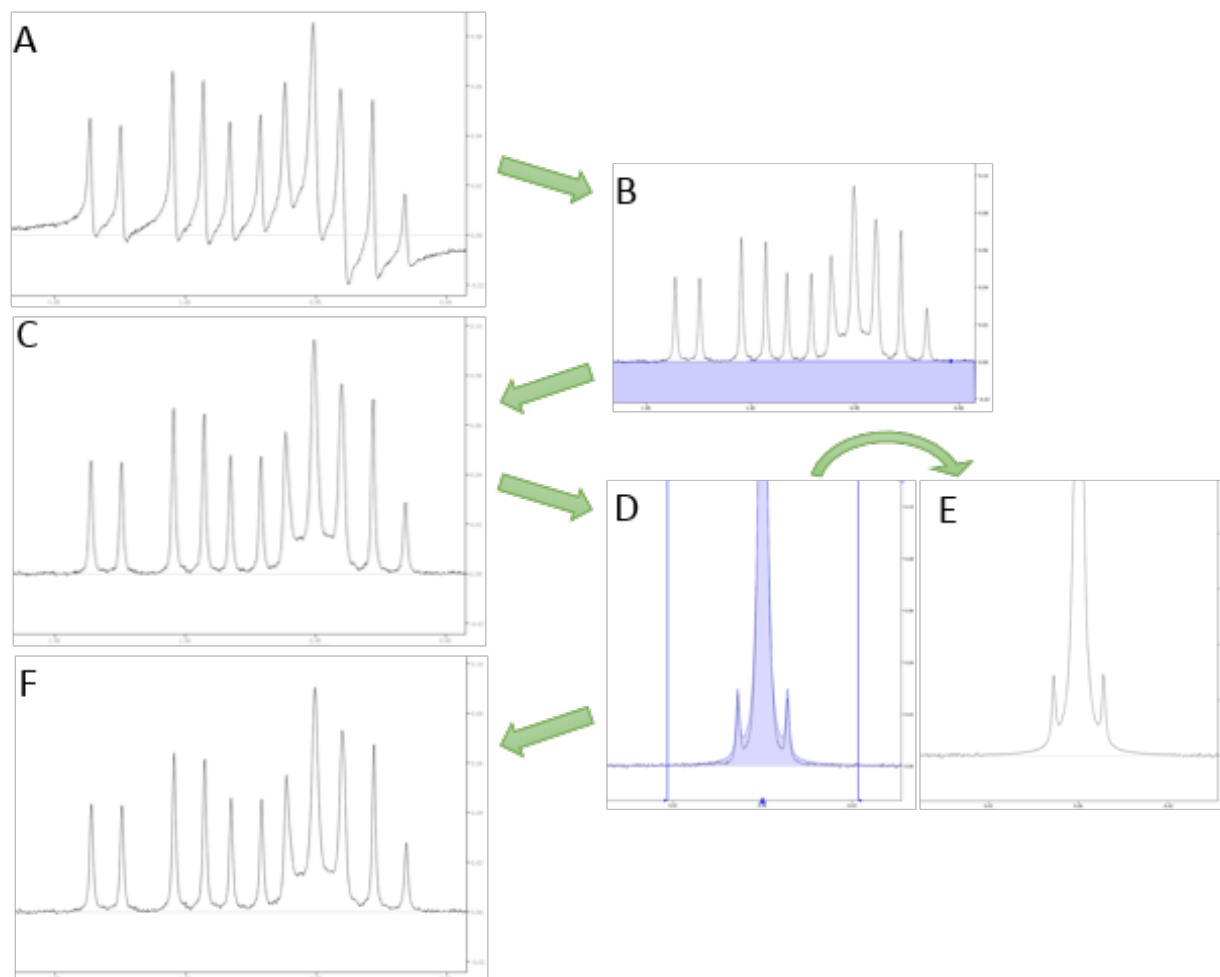


Figure 4-1. Effects of each spectral processing step on the NMR spectrum. The NMR spectrum within the range of 0.9-1.5 ppm of defined media for CHO cells is used as an example. (A) Raw NMR spectrum, no processing steps applied. (B) Applying phase corrections. (C) NMR spectrum after phase correction and baseline correction. (D) and (E) Applying shim correction to DSS signal at 0.0 ppm. (F) NMR spectrum after phase, baseline and shim corrections applied.

4.2.3 Building Metabolite Profiles

The assignment of compound identity to peak signals and quantifying those metabolite concentrations are performed using a targeted profiling technique. This method of interpreting the NMR spectra involves using the templates in the Chenomx NMR suite software library of compounds. By matching the compound templates to spectral signals, the identity of the compound is confirmed and the compound concentration is set (Figure 4-2). Assigning compound identities consistently for a given signal and setting compound concentration with high precision, even in areas of the spectrum with a high degree of signal convolution, are skills that are only developed with practice.

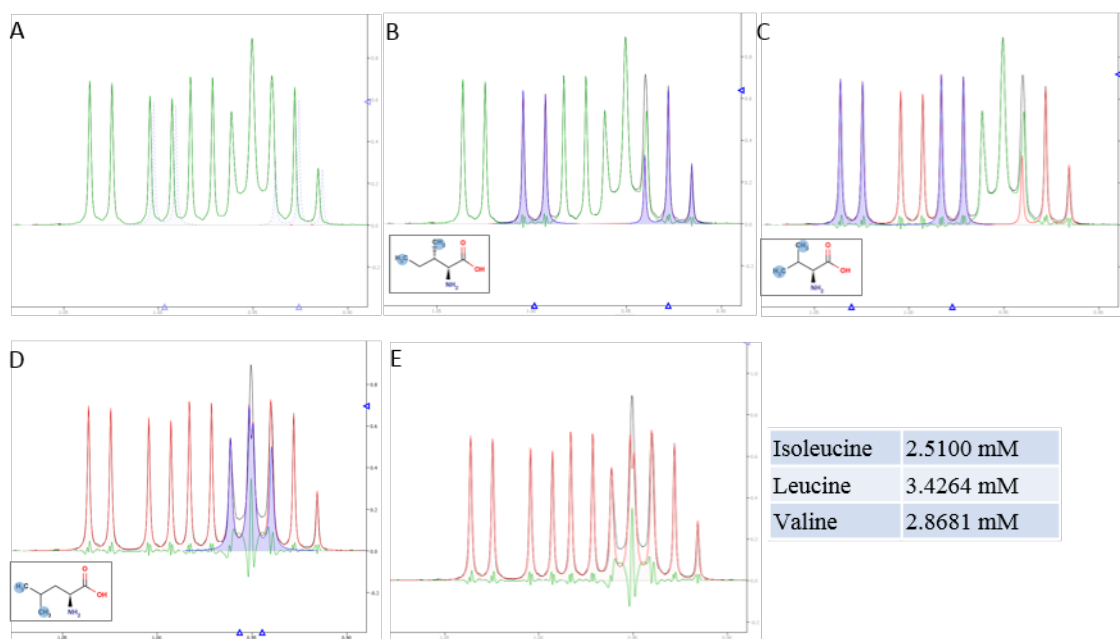


Figure 4-2. Identifying and quantifying compounds in NMR spectrum of defined CHO media (0.9-1.5 ppm). Black line represents the NMR spectrum. Blue dotted line represents the software's compound template, not yet fitted to the spectrum. Compound template fitted to the spectrum to confirm metabolite assignment and concentration to the corresponding signal is shown in transparent blue. Red line represents the sum line, showing the area of the spectrum that has already been fit with a compound, where overlaying templates have an additive effect. Green line represents the subtraction line, which is the difference between the NMR spectrum (black) and the profiled compounds (red). (A) Non-profiled spectrum. (B) Isoleucine. (C) Valine. (D) Leucine. Only clusters between 0.9 and 1.5 for each compound are shown. Signal clusters outside this range would also be used to determine compound identity and concentration.

4.3 Treating Cell Culture Media with UV Irradiation against Adventitious Agents: Minimal Impact on CHO Performance

4.3.1 Journal Article Authorship

Sandi Yen^{1,#}, Stanislav Sokolenko^{1,#}, Bhavik Manocha¹, Ankit Patras², Farnaz Daynouri-Pancino², Eric JM Blondeel¹, Michael Sasges², Marc G Aucoin^{1,*}

¹Waterloo Institute for Nanotechnology, Department of Chemical Engineering, University of Waterloo,

Waterloo, Ontario, Canada, N2L3G1

²Aquafine Corporation, Valencia CA, USA, 91355-4198

#: contributed equally to this work

***corresponding author:** Marc G. Aucoin, Department of Chemical Engineering, University of Waterloo, Waterloo, ON, CANADA N2L 3G1, Tel: 1-519-888-4567 x36084, Fax: 1-519-888-4347, e-mail address: maucoin@uwaterloo.ca

4.3.2 Justification of Original Work

The following journal article was published in Biotechnology Progress on July 29, 2014. The experiment was a collaborative project between the Aucoin group at the University of Waterloo and the Sasges team at Trojan Technologies (London, ON, Canada)/Aquafine Corporation (Valencia, CA, USA). UV irradiation of solutions and related protocols were conducted by Trojan UV, while metabolite analyses and cell culture experiments were performed by the Aucoin group of University of Waterloo. The objective of this experiment was to investigate the effect of UV irradiation on the chemical composition of cell culture media and to determine the level of irradiation that can be

tolerated before negative effects are observed in simple batch growth cell culture. Within the Aucoin group the contributions towards this project must be credited to several team members under the supervision of Marc Aucoin. The cell culture responsibilities were shared between Sandi Yen and Bhavik Manocha. Antibody analyses were performed by Jann Catherine Ang and metabolite profiles of cell culture media were constructed by Sandi Yen. Data analysis and interpretation was performed by Stanislav Sokolenko. The overall manuscript was written primarily by Sandi Yen and Stanislav Sokolenko, with editing provided by Eric Blondeel. This work is included in this thesis because establishing metabolite profiles of unspent defined media was the initial step to being able to build metabolite profiles of complex solutions to be encountered throughout the body of the thesis. Unspent defined media established foundational skills in NMR-based metabolite analysis because it is a solution that contains many metabolites, thus producing complex NMR spectra. At the same time, the known composition of the media provided information on the accuracy and sensitivity of the metabolite profiles. In this respect, the experiment presented here fits appropriately as the preliminary work within this body of research. With minor editorial changes to fulfill formatting requirements, this chapter is substantially as it has been submitted and published in *Biotechnology Progress*.

4.3.3 Abstract

Sterility of cell culture media is an important concern in biotherapeutic processing. In large scale biotherapeutic production, a unit contamination of cell culture media can have costly effects. Ultraviolet (UV) irradiation is a sterilization method effective against bacteria and viruses while being non-thermal and non-adulterating in its mechanism of action. This makes UV irradiation attractive for use in sterilization of cell culture media. The objective of this study was to evaluate the effect of UV irradiation of cell culture media in terms of chemical composition and the ability to grow cell cultures in the treated media. The results showed that UV irradiation of commercial cell culture media at

relevant disinfection doses impacted the chemical composition of the media with respect to several carboxylic acids, and to a minimal extent, amino acids. The cumulative effect of these changes, however, did not negatively influence the ability to culture Chinese Hamster Ovary cells, as evaluated by cell viability, growth rate and protein titre measurements in simple batch growth compared to the cells cultured in control media exposed to visible light.

4.3.4 Introduction

Sterility is an ongoing challenge for commercial bioprocesses, particularly in the production of biotherapeutics, where sterility assurance levels are mandated across entire manufacturing processes. Major pharmaceutical companies have suffered significant losses due to both bacterial and viral contaminations at various stages of production^{93–95}. Such incidents call attention to the need to address the risk profile associated with biopharmaceutical products. As an ingredient in the production process, culture media is a potential source of contamination that can propagate downstream and ultimately affect patient safety. Complementary media sterilization methods can confer additional safeguards for biopharmaceutical production and act as a viral barrier to mitigate contamination risks.

Sterility assurance levels of 6, or a probability of 10^{-6} that a single viable microorganism survives the sterilization process and is present in the final product, is the level of sterility that has become the standard requirement for pharmaceutical products and sterilization processes⁹⁶. This level of sterility is traditionally accomplished via terminal sterilization processes by high heat and pressure^{97,98}, which can be inappropriate for solutions having heat-labile components. High-temperature/short-time (HTST) pasteurization is also widely used as an upstream viral barrier. While this technology has been shown to be effective in inactivating viruses in media^{99,100}, HTST has been shown to be incompatible with certain serum containing media¹⁰⁰. Applying HTST can also cause

phosphate- and calcium-based precipitates to form, ultimately impacting the operation of the HTST system^{101,102}.

As an alternative to heat treatments, aseptic procedures require all production equipment and components to be either 'single-use' disposable or steam-sterilized before use. Cell media are often subjected to sterile filtration with 0.22 μm filters following aseptic steps to further preclude media contamination⁹⁷.

Ultraviolet (UV) radiation is known to be an effective sterilant against various biological contaminants, including bacteria and viruses^{103–106}. For example, active cytomegalovirus (CMV), can be reduced by one log with a UV fluence of 5 mJ/cm^2 ¹⁰⁷, while minute virus of mice (MVM) requires only 2 mJ/cm^2 ^{108,109}. More resistant viruses, such as adenovirus require up to 306 mJ/cm^2 to achieve a log inactivation of approximately 6^{15,16}.

UV disinfection is favourable for its non-thermal and non-adulterating characteristics, and has been adopted by several industries, including the biopharmaceutical industry, for packaging and surface sterilization applications^{112–115}. UV irradiation is able to approach the sterilization criteria of 6-log reduction with MS2 bacteriophage, which is used as a model challenge organism and used to assess reduction equivalent fluence of a UV system^{116–118}. This makes UV irradiation a sterilization process fitting for diverse biopharmaceutical applications, such as in the sterilization of polyethylene bottles as well as the disinfection of air, water and surfaces¹⁰⁴.

The application of UV irradiation to cell culture media as a sterilization method has been limited thus far and conducted primarily in industry. Previous work on UV irradiation of media has suggested that treatment can result in changes to substrate concentrations or the production of various by-products¹¹⁹, possibly leading to cell growth failure¹²⁰ or changes in critical quality attributes

(CQA) of the final product. The objective of this work was to investigate the effect of UV irradiation on the chemical composition of cell culture media and to determine the level of irradiation that can be tolerated before negative effects are observed in simple batch growth cell culture.

4.3.5 Materials and Methods

4.3.5.1 Experimental design

Two culture media were evaluated in this study: a protein-free, serum-free, chemically defined medium optimized for the growth of Chinese hamster ovary (CHO) cells and expression of recombinant proteins in suspension culture, as well as Media 199 (M199), formulation 11150, a fully defined protein-free, serum-free media used to culture chick embryo fibroblasts. Four UV irradiation doses were considered in this work, 110 mJ/cm², 180 mJ/cm², 250 mJ/cm², and 400 mJ/cm². Due to the wide range of UV fluence at which bacteria and viruses are inactivated, these high level dosages were selected to be effective against both bacteria, which are more UV-sensitive, and viruses, which tend to be more UV-resistant¹¹⁰. Fluence was estimated based on the log-reduction of MS2 bacteriophage in the media given that 110 mJ/cm² yields a 5-log reduction of this organism.

To best characterize the overall chemical changes in media, we used a nuclear magnetic resonance (NMR)-based metabolomics approach to characterize a large variety of media metabolites that could potentially be affected by the UV treatment. To achieve a high level of confidence in the observed effect of UV treatment on metabolite concentrations, each treatment was repeated 3-6 times per sample set, with multiple sample sets generated over a 4 month period. This allowed the separation of observation and process variability from the quantification of treatment effect.

The impact of UV irradiation on cell growth was determined by culturing cells in treated media. Two controls were used for the analysis: unmodified media that had not been treated in any

way, which is referred to as samples receiving 0 mJ/cm², and non-irradiated samples that had been stirred and exposed to air for the same duration as the UV-irradiated samples.

4.3.5.2 Media irradiation

Irradiation of media was executed by Trojan Technologies (London, Ontario, Canada). A collimated beam apparatus equipped with a low-pressure mercury lamp emitting at 254 nm was used for uniform and quantified UV irradiation of media¹²¹. The calibration, fluence determination and quality assurance protocols associated with UV irradiating solutions were carried out according to standardized practices in the water disinfection industry¹²². The media was irradiate at low volumes, ~5 mL, in 10-mL beakers with stirring to ensure uniform UV delivery throughout the liquid solution. The actual UV fluence delivered to the media was calculated based on the measured UV dose emitted by the collimated beam and the optical properties of the media being treated.

4.3.5.3 Verification of UV fluence

In order to establish the actual UV fluence values delivered to the culture media, a viral clearance test was conducted using a challenge organism inoculated into the culture media. The challenge organism, MS2, is a single-stranded RNA virus, with an icosahedral shell approximately 27 nm in diameter, used extensively in validation of UV disinfection systems for drinking water¹²³.

MS2 bacteriophage was suspended in phosphate-buffered saline and irradiated in a collimated beam to various UV fluence values. The buffer used in this characterization had high transparency at 254 nm, typically 90% at 1 cm, so that the UV intensity gradient in the fluid was small, and therefore the uncertainty in the UV fluence was small. The 254 nm irradiance at the sample surface was measured using a radiometer with NIST-traceable calibration. The UV fluence gradient caused by optical absorbance of the fluid, along with the effects of beam divergence, non-uniformity, and

surface reflection were all accounted for as recommended in the standard method. The samples were stirred during irradiation to ensure that all organisms were exposed to the same integrated UV fluence over the course of the irradiation. After irradiation, each sample was diluted serially and aliquoted into culture tubes containing 1 mL *Escherichia coli* broth culture and 20 mL of molten tryptone yeast extract glucose agar containing triphenyl tetrazolium chloride. The mixtures were mixed by inversion and plated into sterile Petri dishes. The agar was allowed to solidify and the plates were incubated at 35°C for 18–24 h before performing a plaque assay. The plates were then evaluated to determine the number of plaque forming units. By comparing with the control (non-irradiated) sample, the relationship between UV dose and the log-reduction of this population of MS2 was established.

With the UV-sensitivity of the organism determined, the UV dose applied to the culture media was determined by using MS2 as the dose indicator. This is sometimes denoted as “Reduction Equivalent Dose (Fluence),” since it is inferred from the log reduction of a well-characterized challenge organism. MS2 from the characterized population was inoculated into samples of the culture media, and the optical properties of the inoculated media were measured. The optical properties of the media were used to calculate the irradiation times necessary to achieve a desired UV dose in the collimated beam apparatus. The samples were then irradiated for the prescribed time, and then serially diluted and cultured as described above. The log reduction in numbers of active MS2 were used to calculate the UV dose received by the media in each irradiation by using the sensitivity of the MS2 as established by the buffer tests. The resulting Reduction Equivalent Dose vs. target dose has been plotted and may be seen in Figure 4-3. This relationship was used to determine the UV fluence values used in subsequent irradiations of mammalian cell culture media.

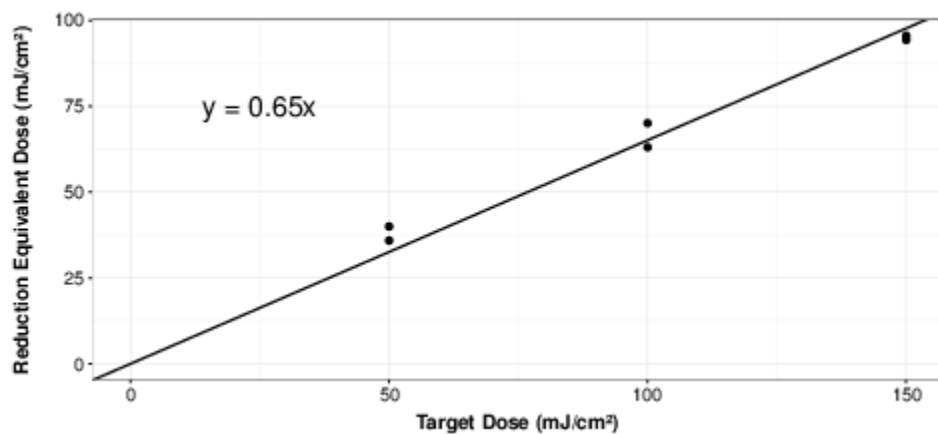


Figure 4-3. Received UV fluence in commercial media measured by bioassay (MS2 bacteriophage) vs. target value. The “Reduction Equivalent Dose (R.E.D.)” applied to the culture media was determined by well-characterized MS2 phage as the dose indicator (Target dose).

4.3.5.4 Cell culture

CHO^{BRI} cells, obtained from the National Research Council Canada, were suspension-adapted cells grown in proprietary, serum-free growth media blend, BioGro-CHO, supplied by BioGro Technologies Inc. (Manitoba, Canada). The cells were adapted to grow in a commercial serum-free, chemically-defined CHO media, supplemented with GlutaMAXTM (Invitrogen Corp., Burlington, Canada) and HT Supplement (Invitrogen Corp., Burlington, Canada). The parental culture was routinely maintained in shaker flasks kept in a humid incubator (37° C, 5% CO₂) and agitated at 100 rpm. Both the parental and experimental cultures were cultured in 125 mL, graduated non-pyrogenic polycarbonate Erlenmeyer flasks (Corning Inc., NY, USA). Once the mother flask reached a viable cell density of 2-3x10⁶ cells/mL, the experimental cultures were inoculated at a seeding density of 0.2x10⁶ viable cells/mL.

For culture tests, the irradiated and control samples of CHO media were supplemented with GlutaMAXTM and HT Supplement after media irradiation and prior to inoculation. GlutaMAXTM and HT Supplement were not subjected to UV irradiation in these experiments. Multiple irradiated samples were pooled to achieve 15-20 mL volumes required for culture growth and supplemented as above.

4.3.5.5 Cell count

Each experimental culture was counted at 24-hour intervals for five days using both the classic haemocytometer counting method, and with a Coulter Counter Z2 (Beckman-Coulter, Miami, USA) for cell size distribution analysis. Samples counted in the haemocytometer were prepared in a 1:2 dilution with Trypan blue (10% w/v in PBS) before loading the diluted sample into the counting chamber. In circumstances where the cell count was greater than 3-4 x 10⁶ viable cells/mL, the sample was diluted by a factor of 10 using CHO media as the diluent, while maintaining a 50:50 ratio

of dye to sample and diluent. Cell size distributions were obtained from a Coulter Counter Z2 by loading samples diluted 100-fold in Isoton® II solution. Cell viability was also calculated daily based on haemocytometer counts.

4.3.5.6 NMR spectroscopy and metabolite profiling

NMR samples were prepared by the combination of 630 µL of media and 70 µL of internal standard composed of 99.9% D₂O with 5 mM 4,4-dimethyl-4-silapentane-1-sulfonic acid (DSS) serving as a chemical shift indicator (CSI) and 0.2% w/v sodium azide to inhibit bacterial growth (Chenomx Inc., Edmonton, Canada). The solution was vortexed and pipetted into 5 mm NMR tubes (NE-UL5-7, New Era Enterprises Inc., Vineland, NJ, USA) for scanning. NMR spectra were obtained with a 600 MHz Bruker Avance spectrometer, equipped with a Triple Resonance Probe (TXI 600). The spectra were acquired using the first increment of the NOESY pulse sequence with a 1-second presaturation pulse, followed by a 4-second acquisition time. All spectra processing was carried out with Chenomx NMR Suite 7.5 (Chenomx Inc., Edmonton, Canada). Baseline, phase, shim, and chemical shift corrections were all performed manually using tools available with the software. Briefly, baseline correction was carried out by the selection of cubic spline points to subtract from the observed spectra, while shim correction corresponded to reference deconvolution to remove line asymmetry¹²⁴. Subsequently, compounds were quantified by targeted profiling. The observed spectra were superimposed with Chenomx's built-in library of chemical resonances, with metabolite concentration quantified using DSS as a reference compound (for more information on targeted profiling, see Weljie et al.³¹).

4.3.5.7 Antibody titer

After 5 days of culturing, the supernatant for the CHO cultures were collected. Each acquired sample was concentrated down to 1 mL volume using a Centriprep Centrifugal filter unit with an Ultracel-10

membrane (EMD Millipore, MA, USA). The concentrated samples were passed through a Protein A HP Spin Trap column (GE Healthcare) and purified IgG antibodies were collected. Antibody titer (total extracellular protein) was determined using a PierceTM BCA protein assay kit (Pierce Biotechnology, Rockford, USA).

4.3.6 Results

4.3.6.1 NMR metabolite analysis

Twenty-eight compounds were identified and quantified in the media by NMR, of which twelve compounds were found to have statistically significant trends with respect to UV irradiation: pyruvate, acetate, arginine, lysine, threonine, tryptophan, tyrosine, choline, methionine, formate, ethanolamine and pyridoxine (in decreasing order of concentration change). When the concentration profiles of these compounds were compared to changes from exposure to visible light and air (without UV) in the control samples, it was found that five compounds had similar changes from visible light and air exposure: arginine, threonine, choline, methionine, and ethanolamine. These changes may be related to exposure to oxygen or to the effects of room light, but are not related to UV exposure, since there was negligible UV exposure in these controls. Concentrations of pyruvate, acetate, and formate were also observed to change with exposure to visible light and air, albeit much less than under UV treatment. Therefore, the observed concentration changes of only seven compounds could be directly linked to UV exposure - pyruvate, acetate, lysine, tryptophan, tyrosine, formate, and pyridoxine. The concentration profiles of these compounds were all approximately linear as a function of UV fluence (Figure 4-4). The specific concentration changes can be seen in Table 4-1, calculated per 100 mJ/cm² of UV fluence. Amino acids changed the least in relative concentration, ranging from decreases of 3% to 7%. A much greater effect was observed for the carboxylic acids, with acetate concentration, for example, increasing twofold.

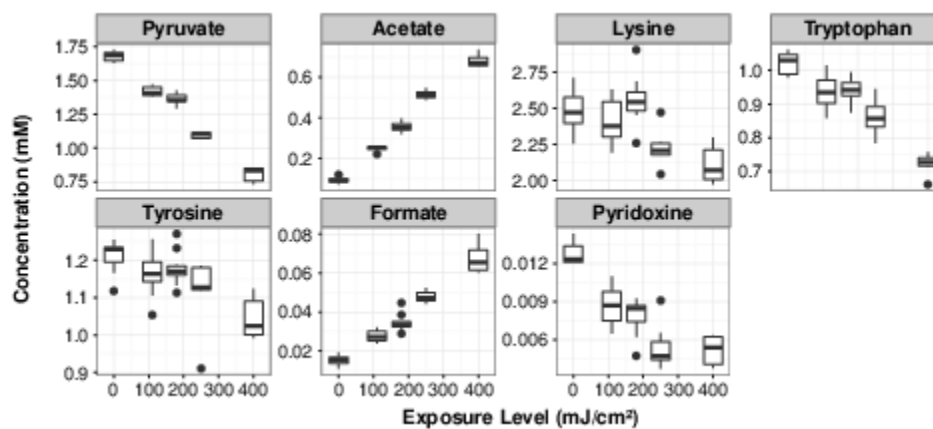


Figure 4-4. Compound concentration as a function of UV fluence. The box-plots contain aggregated data from multiple experiments in which exposures were performed 3-6 times at each level.

Table 4-1. Changes in compound concentration per 100 mJ/cm² of UV dose.²

Compound	Concentration Change (μM)	Concentration Change (%)
Pyruvate	-218.7 +/- 6.6	-13.0 +/- 0.4
Acetate	152.1 +/- 3.4	164.0 +/- 3.6
Lysine	-91.5 +/- 21.4	-3.7 +/- 0.9
Tryptophan	-70.6 +/- 5.6	-6.9 +/- 0.5
Tyrosine	-39.6 +/- 7.2	-3.3 +/- 0.6
Formate	13.2 +/- 0.5	86.5 +/- 3.5
Pyridoxine	-2.0 +/- 0.2	-15.6 +/- 1.6

²Linear regression was performed on the change in compound concentration as a function of UV dose with the presented number corresponding to the predicted concentration change per 100 mJ/cm² (with the +/- representing standard deviation).

4.3.6.2 Cell culture

Cell culture performance was primarily assessed by viability, growth rate, and final cell concentration. Controls (exposed to visible light) were cultured alongside cultures grown in UV-treated media to account for changes in cell behaviour. Overall, UV-treatment of the media did not have a significant effect on cell growth. Therefore, any variability observed in cell growth and protein production can be attributed to random effects imparted by the experimental environment, such as visible light. Cell viabilities of both treated and control samples were routinely above 95%, with no significant difference between the two groups as tested by paired Student's t-test at the 95% confidence level. Additionally, no differences were observed in cell size distributions (data not shown). Growth rates can be seen in Figure 4-4A and the final cell concentrations in Figure 4-5B. Overall, the differences between treated and control cultures were not statistically significant.

In the course of this study, quail muscle fibrosarcoma (QM5) cells were also cultured in M199 media treated with a dose of 110 mJ/cm². No differences were observed between QM5 cells grown in UV-treated and those grown in control media. In both cases, the cells were able to reach confluence within 4 days, with no apparent changes in morphology. Higher UV doses were not assessed for this media.

4.3.6.3 Protein titre

Similar to the growth data, recombinant IgG protein titers from the CHO cells were found to be very similar for UV-treated and control media. Comparing the yield on a purely volumetric basis suggested that UV-treated media may result in somewhat better protein production (Figure 4-5C). However, when integrated viable cell density was used to calculate specific protein production levels on a per-cell basis¹²⁵, there were no significant differences in protein yield (Figure 4-5D).

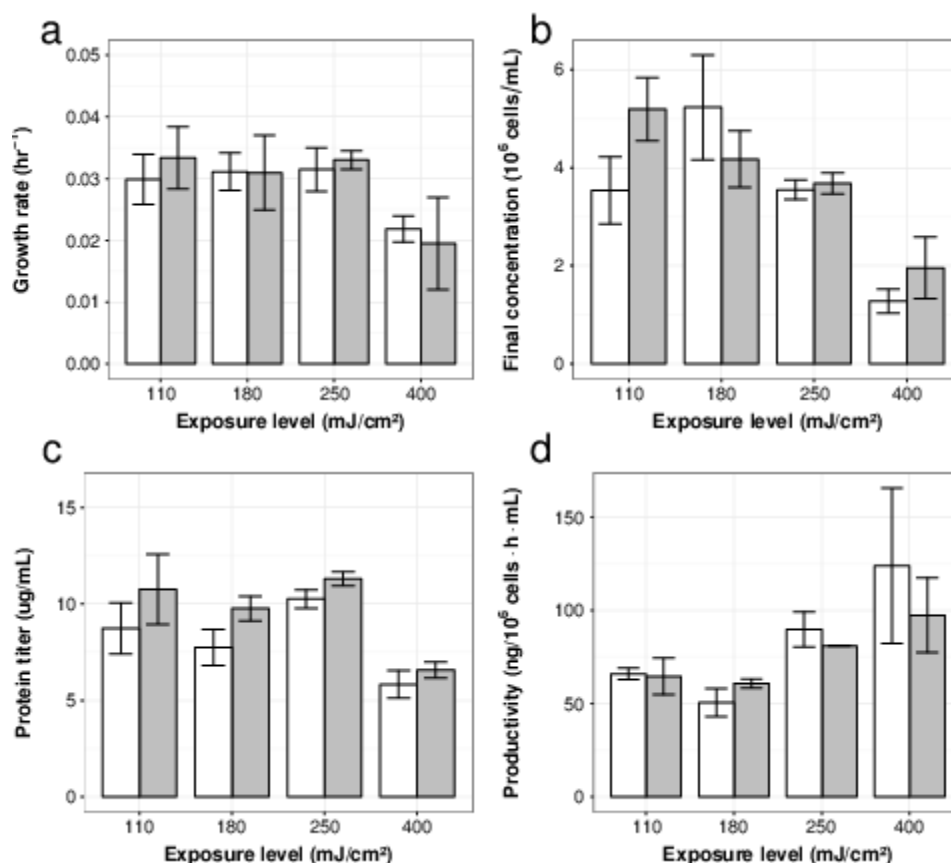


Figure 4-5. (A) Comparison of cell culture performance on control (white) and treated (grey) media. In all cases, controls (exposed to visible light) were cultured alongside cultures grown in UV-treated media to account for changes in cell behaviour over a 4-month period. No statistically significant differences were observed between treated and control media for any of the culture parameters, as assessed with a two-sample t-test at a 95% confidence level. a) Mean exponential phase growth rates (with standard errors) for cell cultures grown in UV-treated and control media. Values were calculated from the log-regression slope of concentration data from days 1-3. (B) Mean final cell concentrations (with standard errors) for cell cultures grown in UV-treated and control media. For the 250 mJ/cm² dose, the cultures grown on treated media resulted in nearly identical final concentrations, leading to overlapping error bars. (C) Mean protein yield (with standard errors) for cell cultures grown in UV-treated and control media. (D) Mean specific protein yield (with standard errors) for cell cultures grown in UV-treated and control media. Values were calculated by dividing the volumetric protein yield by the integrated viable cell density. For the 250 mJ/cm² dose, the cultures grown on treated media resulted in nearly identical specific protein yield, leading to overlapping error bars.

4.3.7 Discussion

Although UV irradiation has been previously reported to negatively impact cell growth¹²⁰, the limiting UV dose values were not indicated. No negative effect was observed in this work for either CHO cells grown in UV-irradiated CHO medium or QM5 cells grown in UV-irradiated M199, suggesting that the impact of UV-irradiation cannot be generalized, or the effect is only encountered when pushing the culture beyond a simple single batch process. The metabolic changes observed due to UV irradiation of media did not impact the overall pH measurement of the media solutions used. It should be noted that we have observed pH shifts in other irradiated media (data not shown), though none in this work. This is in part due to the sodium bicarbonate buffer in the media and the relatively small absolute changes of impacted compounds, such as acetate.

Observed changes in substrate concentrations were also limited. Amino acid concentration changes that did occur were on the order of 1-7% per 100 mJ/cm². Most cell culture media contain these compounds in excess so the extent of degradation of amino acids from UV irradiation, even at high UV doses, is unlikely to impact culture performance. The small changes in substrate concentration also meant that any specific degradation products were generally masked by overall compound variability, which we have previously found to be approximately 5%-10% for most compounds, measured as a coefficient of variation³⁹. Pyridoxine was the only vitamin that could be observed with NMR, and therefore no general conclusions can be made about overall vitamin stability. However, the cell growth data indicates that the level of vitamin degradation was too low to impact cell growth compared to our control.

The largest changes in compound concentration in these media were observed for carboxylic acids. Pyruvate concentrations were found to decrease by approximately 13% per 100 mJ/cm², with acetate and formate increasing by 160% and 90% per 100 mJ/cm², respectively. This corresponded to

absolute changes of 218 μM , 152 μM , and 13 μM for pyruvate, acetate, and formate respectively. The three compound concentrations showed parallel trends in both control (visible-light-exposed) and treated media. Considering the high degree of precision in its quantification (2% coefficient of variance including set-to-set variability), pyruvate concentration can be seen as an attractive marker for assessing the overall effect of UV treatment in cell culture media using NMR spectroscopy. As previously mentioned, no additional impact on cell growth was observed compared to the control in spite of such changes to these metabolites.

This decrease in pyruvate is not altogether surprising. Pyruvate has been previously identified to be a prominent antioxidant in biological systems and has been found to yield acetate when reacting with peroxynitrite, a powerful oxidant¹²⁶, or hydrogen peroxide¹²⁷. It may serve a similar role in the quenching of oxidative products formed during UV treatment. Based on absolute concentration changes, it would appear that acetate is the major product of such reactions, although formate may also be involved.

The UV doses applied in this study would achieve disinfection for most adventitious agents relevant to CHO cell culture, such as mouse minute virus and reovirus^{95,103,110}. Since filtration is seen as an effective method to remove larger virus and bacteria, UV disinfection is primarily seen as a complementary method to inactivate adventitious agents not captured by filtration.

4.3.8 Conclusion

The UV doses applied in this study did not functionally impair the culture media tested, as measured by cell growth and antibody production of CHO cells. Cell morphology was not impacted and the cell size distribution was comparable between cultures grown in UV-treated media and control media that had only been exposed to visible light. These findings were consistent at various fluence levels, even

at very high exposure levels of 400 mJ/cm². QM5 cells were also unaffected by media irradiation, though media was only exposed to UV doses up to 110 mJ/cm². Changes to major media components were largely limited to carboxylic acids, with only minimal changes in amino acid concentrations. The findings in this study serve as a demonstration of some of the changes that can occur to cell culture media after applying UV sterilization. While these results cannot be generalized to all cell lines and culture media, this study provides insight into the effects of UV-sterilization in cell culture applications. Further validation with investigations focusing on longer culture duration, higher cell densities, repeated passaging in UV-treated media, and ultimately the effect on protein product critical quality attributes is currently being undertaken.

4.3.9 Acknowledgements

This work was supported in part by NSERC ENGAGE and NSERC Strategic Network (MabNet) Grants to MGA and NSERC Canada Graduate Scholarship to SS. The authors also thank Jann Catherine Ang for quantifying the antibodies produced in this study.

Chapter 5 Standardization of Fecal Water Samples and Analysis

5.1 Chapter Objective

The focus of this chapter is on setting a level of standardization on reporting on the fecal metabolome. Fecal water is the aqueous phase of fecal samples that is prepared by adding water (buffer) to feces, mixing and recovering the liquid phase after removing the solids and centrifugation. Even though fecal water is the accepted means of studying gut microbiota²⁻⁴, there is still a need to developing a standard means of preparing and reporting on the biofluid for metabolomic purposes. Granted, there have been investigations into various approaches for studying fecal water using other metabolomic tools^{20,21}. To develop an appropriate standard for reporting on fecal metabolome, metabolite profiles of fecal water obtained from varying water to feces ratio were examined to determine the minimal concentration of feces required for metabolite analysis of fecal water, and propose normalization techniques for reporting on fecal metabolomes.

5.1.1 Materials and Methods

5.1.1.1 Preparation of Fecal Water Samples

Fecal water samples were prepared as follows. Stool samples were collected from Donor E, a neurotypical, physiologically healthy female (Table 3-3) with no dietary restrictions or any use of antibiotics within 6 months of sample collection. After collection, samples were stored at -80 °C. Frozen stool samples were thawed and the desired amount was weighed into a bag, along with autoclaved PBS. The stool-PBS mixture was homogenized in a stomacher for 1 minute, and the slurry was subsequently ultracentrifuged for 2 hours at 50 000 rpm, 4°C. The supernatant was extracted by needle and syringe. The supernatant in the syringe was shaken vigorously and then passed through 0.22 µm syringe filter. The samples were prepared in increasing concentrations of feces, between the ranges of 8-35% (w/v) (Table 5-1), and aliquoted into duplicate samples. Fecal water samples were stored at 4°C until transport (on ice) to the Aucoin lab the next day. The samples were stored at -80 °C upon receipt of the samples until NMR analysis 12 days later. On the day of NMR scanning, fecal water samples were thawed at room temperature and prepared for NMR spectroscopy as outlined in the Chapter 3.

Table 5-1. Fecal concentration of fecal water samples. Duplicate samples were provided for each fecal concentration.

Mass of raw stool samples (g)	Volume of PBS (mL)	Fecal concentration (g/mL)	Fecal concentration (w/v %)
2.6	30	0.09	9%
5.6	30	0.19	19%
15.2	60	0.25	25%
8.2	30	0.27	27%
9.1	30	0.30	30%
20.6	60	0.34	34%

5.1.1.2 Statistical Analysis

Trends in metabolite concentration due to changes in fecal concentration were evaluated for linearity by using linear regression. Correlations were evaluated with Pearson product-moment correlation coefficient, r . Normalizing compound concentrations to be independent of fecal concentration was calculated by dividing metabolite concentration by fecal concentration. Normalizing compound concentrations to acetate, an internal metabolite was calculated by dividing the metabolite concentrations by the concentration of acetate from the same metabolite profile.

5.1.2 Results

5.1.2.1 Fecal Water Analysis and Sample Preparation

In total, 39 compounds were identified and quantified in each of the metabolite profiles of the fecal water samples. An overview of the profiles indicated that there was a linear relationship ($p \leq 0.05$ at the 95% confidence level) between the fecal concentration and the metabolite concentrations. Furthermore, there was a strong positive correlation (correlation coefficient > 0.74) between the fecal concentration at which the samples were prepared and the observed compound concentrations (Figure 5-1). The only compound that did not follow this trend was ethanol. Ethanol did not have a significant linear relationship between compound concentration and fecal concentration. The changes in compound concentrations were apparent even by comparing the NMR spectra directly (Figure A1).

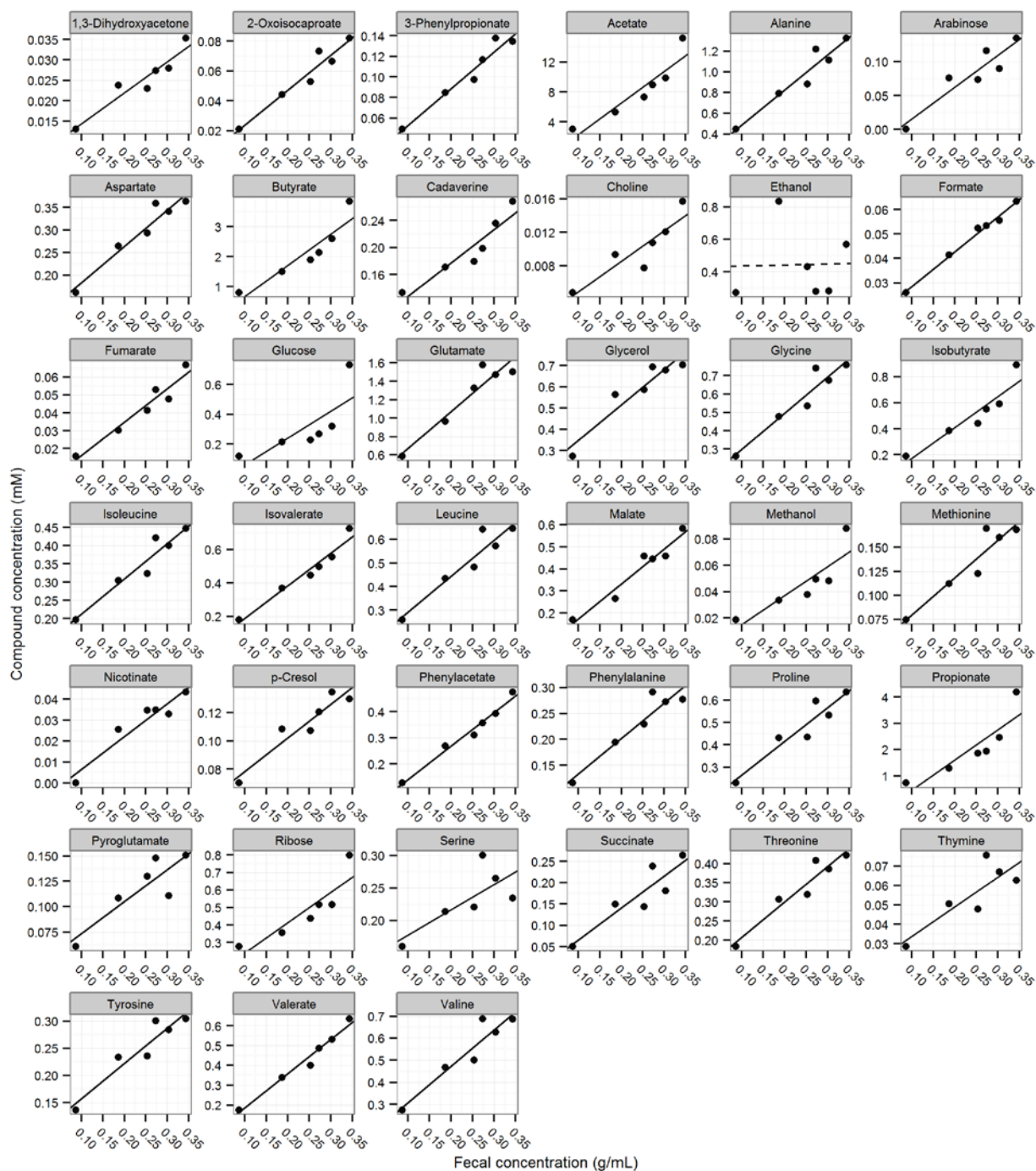


Figure 5-1. Metabolite profile overview. Compounds exhibiting significant linear relationship and a strong correlation between metabolite concentration and fecal concentration are identified by illustrating the trends in a solid black line. Otherwise, a dashed line is used. Mean concentration from duplicate samples shown.

5.1.2.2 Normalization of Metabolite Concentrations

The linear relationship between metabolite concentration and fecal concentration means that metabolite concentrations can be normalized as metabolite concentrations per unit of feces. Dividing the metabolite concentration by fecal concentration expresses the metabolites in terms of mmol/g of feces. Given that all the samples examined here were identical in composition and different only in their dilution, the normalized metabolite concentration should be very similar across different fecal concentrations. Indeed, the spread of the normalized values is quite narrow for each compound (Figure 5-2A).

Alternatively, metabolite concentrations can be normalized to an internal compound such as acetate. Since acetate was a compound that was consistently found at high levels (~4 times larger than the next most abundant compounds), it is an attractive marker to gage the contents of the fecal samples. Normalizing the metabolite concentrations to acetate was performed by dividing each metabolite concentration by the mean acetate value for each fecal concentration. Reporting compounds as a percent of acetate presented the metabolite concentrations independent of fecal concentration. The success of acetate-normalization was demonstrated in the narrow range observed in the acetate-normalized values for every compound. (Figure 5-2B).

In comparing the efficacy of each type of normalization, it was noticeable that the spread of normalized values, as measured by range, in acetate-normalized values is generally smaller than that of the normalization by fecal concentration. The reason that range is used for the comparison instead of standard deviation was because standard deviation is a measurement that is magnitude dependent. Inherent to the normalization method, acetate-normalized values must be between 0 and 1, while the fecal concentration-normalized values must be within 1-3 orders of magnitude of the original mM

concentration. The scale differences between the two types of normalization limit the meaningfulness of comparing the standard deviation values directly.

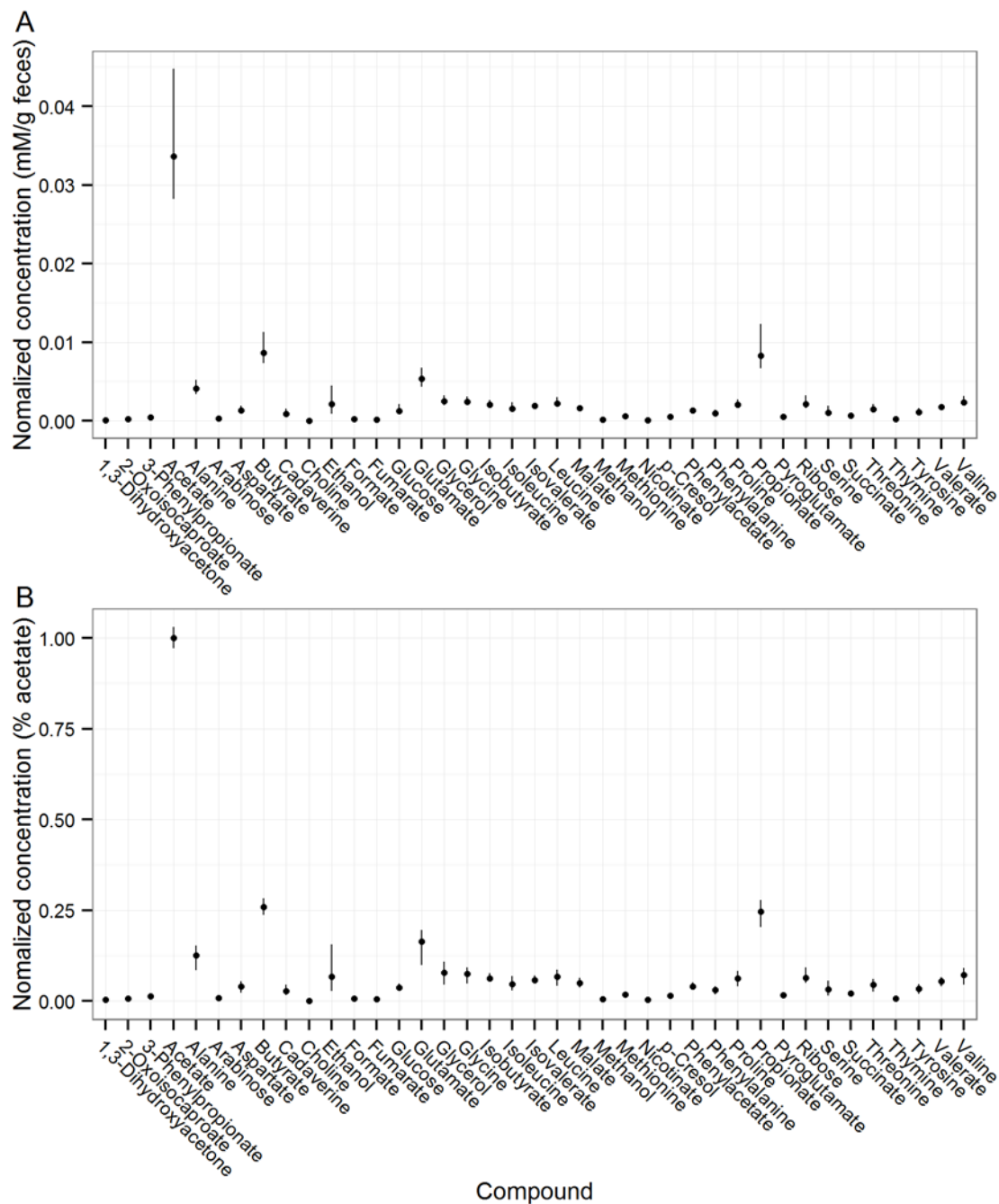


Figure 5-2. Normalized compound concentrations. Bars represent range of normalized values. (A) Normalized by fecal concentration. Calculated by dividing observed compound concentration (mM) by fecal concentration (g/mL). (B) Normalized by mean acetate concentration. Calculated by dividing observed compound concentration (mM) by mean acetate concentration for each fecal concentration.

5.1.2.3 Impact of NMR Spectra on Metabolite Profiles

As mentioned before, 39 compounds were included in the metabolite profiles of the fecal water samples. Even though all 39 compounds were assigned and quantified with relatively high confidence, certain compounds became increasingly difficult to profile as fecal concentration decreased. The decreased resolution, related to the signal-to-noise ratio of the spectra, can yield larger variance in compound quantification, or even impair proper assignment of compound identification. Examples of compounds that had signals that were sensitive to low fecal concentrations included proline, choline, pyroglutamate, p-cresol, arabinose, and glucose (Figure 5-3, Figure A2). For example, proline and glucose both had clusters within the 3.2-3.55 ppm range of the spectrum that were crucial for both identification and quantification of these two compounds, as they are relatively free of convolution with other metabolite signals in this region. As seen in Figure 5-3, proline and glucose clusters were far less prominent relative to the spectrum's noise levels in the low fecal concentration sample compared to the high fecal concentration sample. So, in the spectra representing higher concentrations of feces, compounds such as proline and glucose could be assigned and quantified with relatively high accuracy and confidence. However, if the low fecal concentration samples were profiled independently, i.e. in the absence of also profiling high fecal concentration samples, these compounds could be lost due to the lack of evidence for profiling those subtle signals.

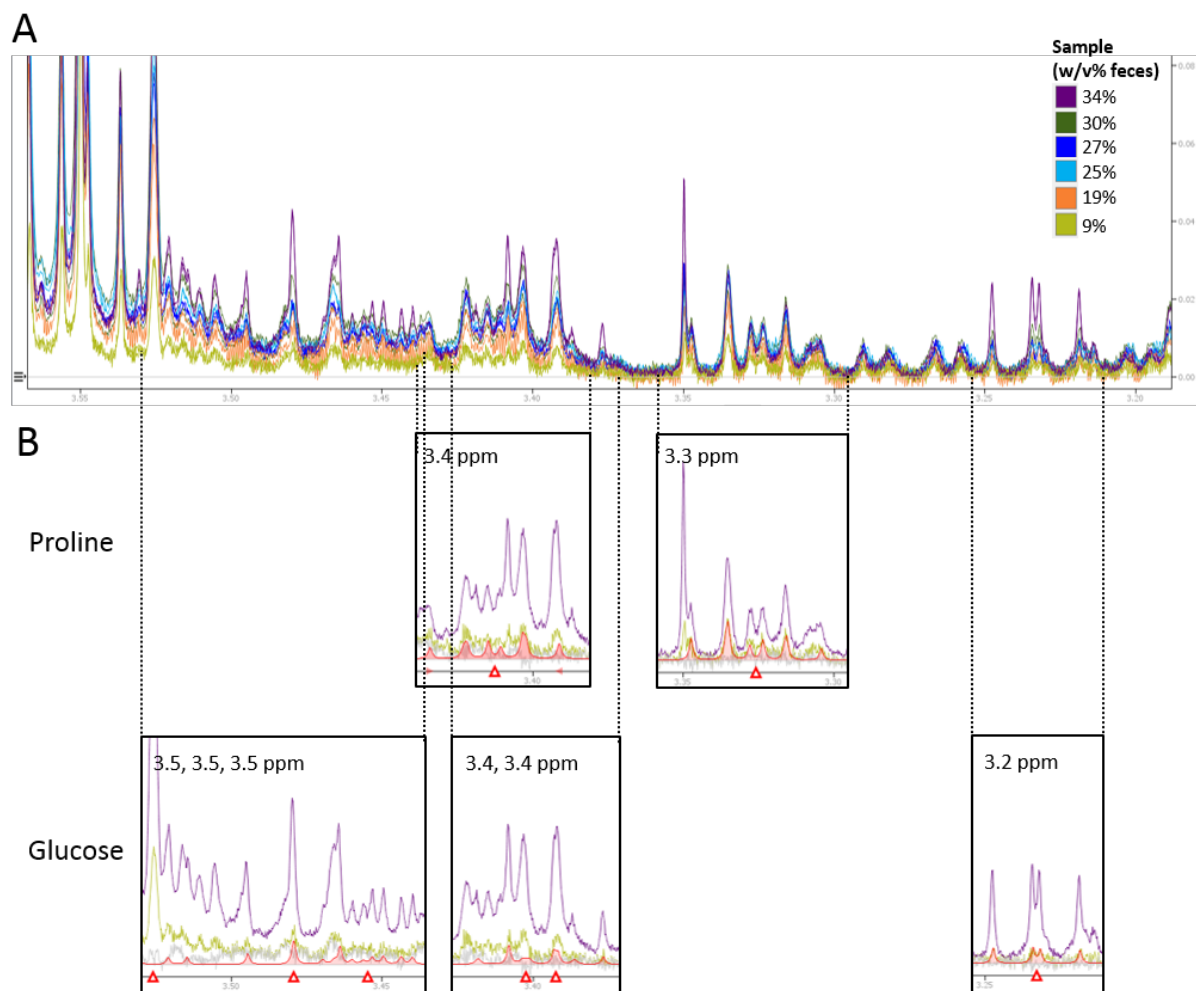


Figure 5-3. Profiling confidence impacted for certain compounds. (A) Spectrum overlay of metabolite profiles of samples for different fecal concentration (3.2-3.55 ppm). (B) Comparison of spectra for lowest fecal concentration (9%) with highest fecal concentration (34%); specifically looking at signals for proline and glucose – two compounds which had profiling confidence sensitive to fecal concentration. Metabolite signals highlighted in red represent concentrations profiled to fit the low fecal concentration spectrum. Light grey line represents the subtraction line between the 9% fecal water spectrum and all profiled metabolites. Note: only the proline and glucose NMR signal clusters found within 3.2-3.5 ppm are presented here.

5.1.3 Discussion

5.1.3.1 Trends in Metabolite Profiles

The metabolite profiles followed a linear trend as fecal concentrations decreased, as expected. Ethanol was the only compound that did not have a strong, linear relationship between compound concentration and fecal concentration. Currently the reason for such a randomized pattern of ethanol concentration observed is unclear. While vapour pressure could be a contributing factor, other volatile compounds, such as methanol, which has a higher vapour pressure than ethanol¹²⁸, did not exhibit random patterns.

5.1.3.2 Normalization of Metabolite Concentrations

The metabolite concentrations were normalized to the amount of feces present in the samples by dividing the metabolite concentrations by the concentration of feces used to prepare the samples. The normalized compound concentrations were calculated to be very similar across samples of various fecal concentrations, re-enforcing the linear relationship between the two parameters. While normalizing the metabolite compounds by fecal concentration enables one to compare samples prepared from the same stool sample, it does not address the variability in water content from one stool sample to the next. Not accounting for the water content of the original stool sample leaves dilution as a factor that is not fully controlled in preparation of fecal water samples, and could ultimately skew the metabolite measurements. One means of overcoming this challenge is to estimate the water content of feces by weighing the stool sample before and after freeze drying¹²⁹. Therefore, a true means of normalizing metabolite measurements to fecal matter content would account for dilutions made during the fecal water preparation process as well as dilution due to the water content of the original stool sample.

Alternatively, reporting the metabolite concentrations independent of feces content circumvents measurement of water content and fecal water dilutions. Normalizing the metabolite profile to the most prominent compound, acetate, presented the metabolites independent of fecal concentrations, while still able to make comparisons across several fecal donations. Acetate is an attractive normalization marker as it is one of the main microbial fermentation products and has been consistently reported as one of the most abundant compound found in fecal metabolomes^{3,61,130}. Even in *in vitro* conditions designed to favour production of other SCFAs, acetate is still present in a relatively significant amount¹³¹. In performing this normalization, the acetate-normalized compound concentrations had small standard deviations, which corresponded to the linear regression model for compound concentration vs fecal concentration. In fact, compared to normalization by fecal concentration, acetate-normalization yielded a narrower spread of normalized values obtained from samples varying in amount of feces content. It should be noted that since acetate-normalization expresses compounds as a ratio, evaluations on acetate levels and comparisons of metabolite levels based on magnitudes are masked. Therefore, it would be ideal to report on both the compound concentrations, as well as their ratio values.

5.1.3.3 Impact of Spectral Limitations

The “profile-ability” of the NMR spectra of fecal water samples should also be taken into consideration when choosing the appropriate fecal concentration. As the fecal concentration decreases, the signal-to-noise ratio of low-concentration compounds worsens and the accuracy and confidence at which the compound can be quantified and identified deteriorates. Proline, choline, pyroglutamate, p-cresol, arabinose and glucose were highlighted as compounds with which the profiling confidence was sensitive to low fecal concentrations. Even though they were not necessarily the compounds found at the lowest concentration in the profiles, they had lower thresholds for

limitations of quantification and identification due to their chemical shift, surrounding peak signals, and signal-to-noise ratios of peaks. Essentially, at low fecal concentrations, the metabolite concentrations of certain compounds were no longer within the limit of quantification, and in the case of arabinose, within the limit of detection. These profiling concerns pertain mostly to the fecal water samples of 9% and 19% fecal concentration. In comparing the samples containing fecal water within the range of 25%-34% fecal concentration, the higher fecal concentration samples did not differ greatly from the mid-range fecal concentration, in terms of ability to distinguish compound signals from noise signals.

5.1.3.4 Minimum Level of Feces Required in Fecal Water Samples

As mentioned earlier, the composition of fecal samples can be highly variable; for example, a gram of one stool donation could have more water content than one gram of another stool donation. Instead of evaluating a fecal water sample based on its concentration of feces, the sample could be identified based on an internal compound. Acetate is a compound that has been consistently observed at high concentrations in fecal water samples, and could be used not only to normalize the compound concentrations, but also to gage the contents of the fecal water sample. The samples of decreasing fecal concentration demonstrated that fecal concentrations at 25% and higher had metabolite profiles that sufficiently cleared the limit of quantification and limit of detection for all the compounds. This minimal fecal concentration could also be expressed as 7.3 mM of acetate, or 0.43 g/L of acetate. The linear relationship between fecal concentration and acetate concentration as well as the ubiquity of acetate in fecal water sample enables this compound to act as a proxy in evaluating the contents of fecal water samples.

Chapter 6 Metabolomic Characterization of Gut-Mimicking CSTR

6.1 Chapter Objective

The purpose of this chapter is to use a metabolomic approach to evaluate the *in vitro* tools used to study the human gut microbiota. Two experiments were performed with the aim of establishing further methodological understanding of both fecal water samples and community cultures grown *in vitro* in an anaerobic continuous stirred tank reactor (CSTR). The first experiment focused on developing a protocol for preparing samples of gut microbiota grown *in vitro* for NMR-based metabolomic analysis purposes. In determining the effects of different stages of sample preparation, such as syringe filtration, centrifugation, etc., a protocol for NMR-based metabolomic analysis was developed. The second experiment presented in this chapter addresses the limited degree standardization of the fecal inocula used for *in vitro* community cultures. Here, the impact of the inoculum on the resulting *in vitro* culture was investigated by culturing several defined communities derived from inocula that contained the same bacterial species, but varied in the proportions at which the bacteria they were present in the inoculum.

6.2 Gut-mimicking CSTR Effluent Sample Preparation and Analysis

6.2.1 Experiment Objective

Samples collected from the gut-mimicking bioreactors require an initial processing step to remove cell material and large media particulates before being suitable for metabolite analysis. The purpose of the following experiment was to determine the effects of various preparation protocols on the metabolite profiles. Two techniques – ultracentrifugation or syringe filtration – were examined to establish an appropriate sample preparation protocol.

6.2.2 Materials and Methods

6.2.2.1 Experiment Design

Along with testing five different ultracentrifugation speeds, two filtration procedures were also tested to see if they could be used as an alternative method for cell and particulate separation. The experiment was conducted to evaluate the effects of various stages in the sample processing procedure. In a first step, some samples were subjected to a “soft spin” (4000 rpm 4, °C, 3 min) to remove the bulk of large insoluble material (pre-processing stage) before being subjected to higher centrifugal forces. In a second step (processing stage), samples were subjected to either ultracentrifugation (4 °C, 2 h) or syringe filtration (35 mm). Following these two stages, all samples were subjected to a post-processing stage consisting of filtration, through a 0.22 µm syringe filter. The different conditions examined are outlined in Table 6-1.

Table 6-1. Sample treatment conditions.

Pre-Process		Process		Number of Replicates	
Soft spin		Ultracentrifuge		Experiment 1	Experiment 2
		rpm	g		
Yes		10 000	7 360	3	3
Yes		20 000	29 430	3	3
Yes		30 000	66 230	3	3
Yes		40 000	117 730	3	3
Yes		50 000	183 960	3	3
No		10 000	7 360	3	0
No		50 000	183 960	3	0
Soft spin	Syringe filtration				
Yes	Serial filtration (1.0µm, 0.80µm, 0.45µm)			3	0
Yes	Single-step filtration			3	8

6.2.2.2 Sample Collection and Handling

An effluent reservoir was collected and stored at 4 °C on the 28th day post-inoculation from a CSTR system propagating feces-derived cultures obtained from Donor A. This reservoir was mixed with a stir bar and sampled in 30-mL volumes for subsequent processing. The effluent reservoir ensured that variability seen in the metabolite results originated from processing differences, and not due to sample-to-sample variability. Ultracentrifugation was performed on Beckman L8-55M, with a Beckman Ti70 rotor and 30-mL Beckman polyallomer Optiseal tubes. Once processed, samples were scanned by NMR.

6.2.2.3 Statistical Analysis

The precision of the metabolite profiles was assessed by analyzing the relative standard deviation (RSD) of the control, non-ultracentrifuged samples in Experiment 2 (Appendix B1.). The two types of filtration procedure, serial and single-step filtration, were compared to each other using Student's t test at the 95% confidence level. A linear regression was used to model this data from all the processing methods (serial filtration, single-step filtration and ultracentrifugation), and the significance of this model was tested by ANOVA. Finally, the metabolite profiles from both experiments were modelled to determine which factors affected the trends observed. This was done using a multiple linear regression (MLR) model using the equation shown in Appendix B2. Since the majority of the compounds demonstrated linear trends, the non-linear terms in this model were not evaluated. Furthermore, the experiment-process interaction term was not tested since none of the samples in Experiment 2 were subjected to the soft spin pre-processing step. The three-factor interaction term was also not evaluated since the small number of replicates would be unlikely to yield enough information to assess such an interaction. The significance of each term in this complex model was evaluated at $p \leq 0.0009$, which is the alpha value of 0.05 after it has been subjected to the

Bonferroni correction (Appendix B3.) to reduce sensitivity to type I error. In this way, a more simplified model was determined. The simplified model produced by each iteration of model simplification was compared to the previous model by ANOVA to assess whether the two models were statistically different. The simplified model comments on the effects of ultracentrifuge speed, experiment, and the pre-processing step on the metabolite profiles. The effects of ultracentrifuge speed was evaluated using the percent change in concentration per 10 000 rpm for each metabolite. This was obtained by using the slope and predicted concentrations from the simplest MLR model (Appendix B4.). The absolute percent decrease in concentration per 10 000 rpm was also analyzed as a function of molecular weight of the corresponding compound. The importance of bias effects was determined as the absolute percentage of the bias coefficient divided by the intercept of the simplified MLR model.

6.2.3 Results

6.2.3.1 Profiling Variability

The samples subjected to single-step syringe filtration in Table 6-1 were used to assess the precision of the metabolite. Since the eight replicate samples were aliquots of the one sample that had been subjected to the preparation protocol, the only source of variability would be the steps involved in constructing the metabolite profiles. The precision of each compound was evaluated based on its relative standard deviation, where the acceptable level of precision was set at 15%, or if near the limit of quantification, 20%. The only compounds that had a RSD value greater than 15% were p-cresol (20.6%) and thymine (16.7%) (Figure A4). Indeed, both p-cresol and thymine were near their limits of quantification in the samples assessed for precision.

6.2.3.2 Processing Methods

Ultracentrifugation speeds, serial filtration (1 μm , 0.8 μm , 0.45 μm and 0.22 μm) and single-step filtration (0.22 μm only) were tested. One of the samples subjected to single-step filtration in Experiment 2 had an unusually high concentration of formate (>0.3 mM), whereas all concentrations of formate in other filtered samples, single-step or serial, did not exceed 0.13 mM. Since this point was an outlier, it was excluded from all analyses. Between the serial filtration and single-step filtration processing methods, no compounds were found to be statistically different. Since the two filtration process methods did not create significantly different metabolite profiles, they were grouped as one set of “control” samples (0 rpm) for the ultracentrifuge-based processing methods.

6.2.3.3 Sources of Effects Impacting Metabolite Concentration

The metabolite profiles were modelled to describe the observed trends in compound concentrations. This model was formed by eliminating non-significant terms from a complex, general model until the simplest model that described the relationships remained. As seen in Figure A3, no compounds in the two interaction effects and the pre-processing effect were significant, revealing that these factors did not contribute to the concentration trends observed, and thus should be removed from the model. The simplest model illustrated that the trends observed were the effects of ultracentrifuge speed and experiment only. This indicated that the effect of the ultracentrifuge speed, which was significant for almost every compound, was independent of other factors, such as experiment and preprocessing steps. Similarly, the effect of experiment, where it was significant, was also independent of all other explanatory factors.

6.2.3.4 Effect of Ultracentrifuge Speed

The simplest MLR model that described the metabolite data showed that almost all compounds in the metabolite profiles demonstrated significant linear relationships with ultracentrifuge speed. These compounds demonstrated a negative linear relationship between compound concentration and ultracentrifuge speed, except for p-cresol, which had a positive linear relationship in both experiments (Figure 6-1). The compounds that did not exhibit significant linear changes with increasing centrifugal force are shown in Table A4. While the linear relationships exhibited by most of the compounds were statistically significant, the percent decrease in metabolite concentration per 10 000 rpm was on average $1.8\% \pm 0.87\%$ in Experiment 1 and $1.8\% + 0.90\%$ in Experiment 2. These linear decreases in concentration with ultracentrifuge speed appeared to be directly related to the molecular weight of the corresponding compound in Experiment 1. However, this effect was not reproduced in Experiment 2, where no significant relationship was observed between the absolute percent decrease and molecular weight (Figure 6-2). The only compound to increase in concentration with centrifugal force was p-cresol, which exhibited a 49.5% increase per 10 000 rpm in Experiment 1 and 41.3% increase per 10 000 rpm in Experiment 2.

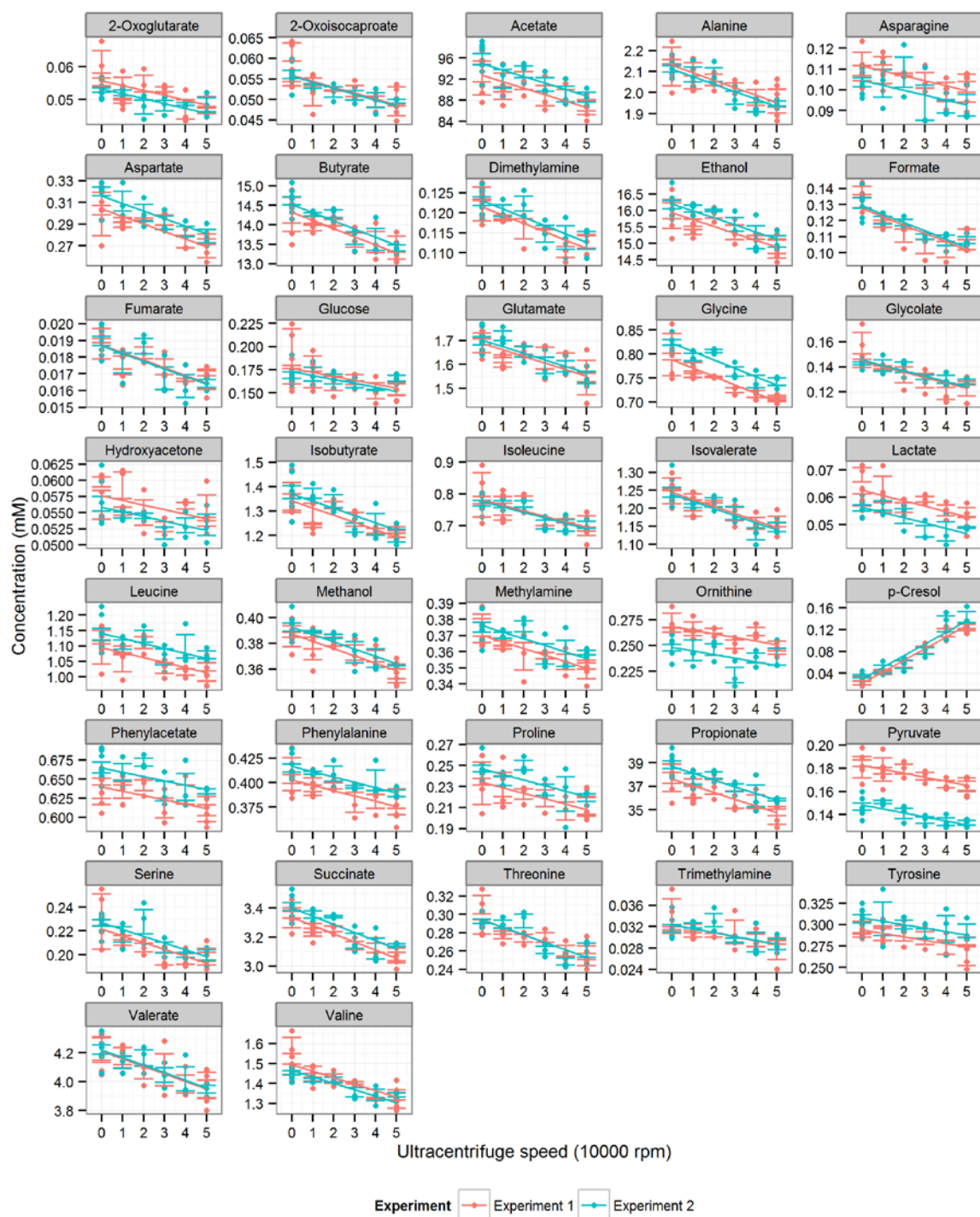


Figure 6-1. Compounds that demonstrate significant linear relationship with centrifugal force.

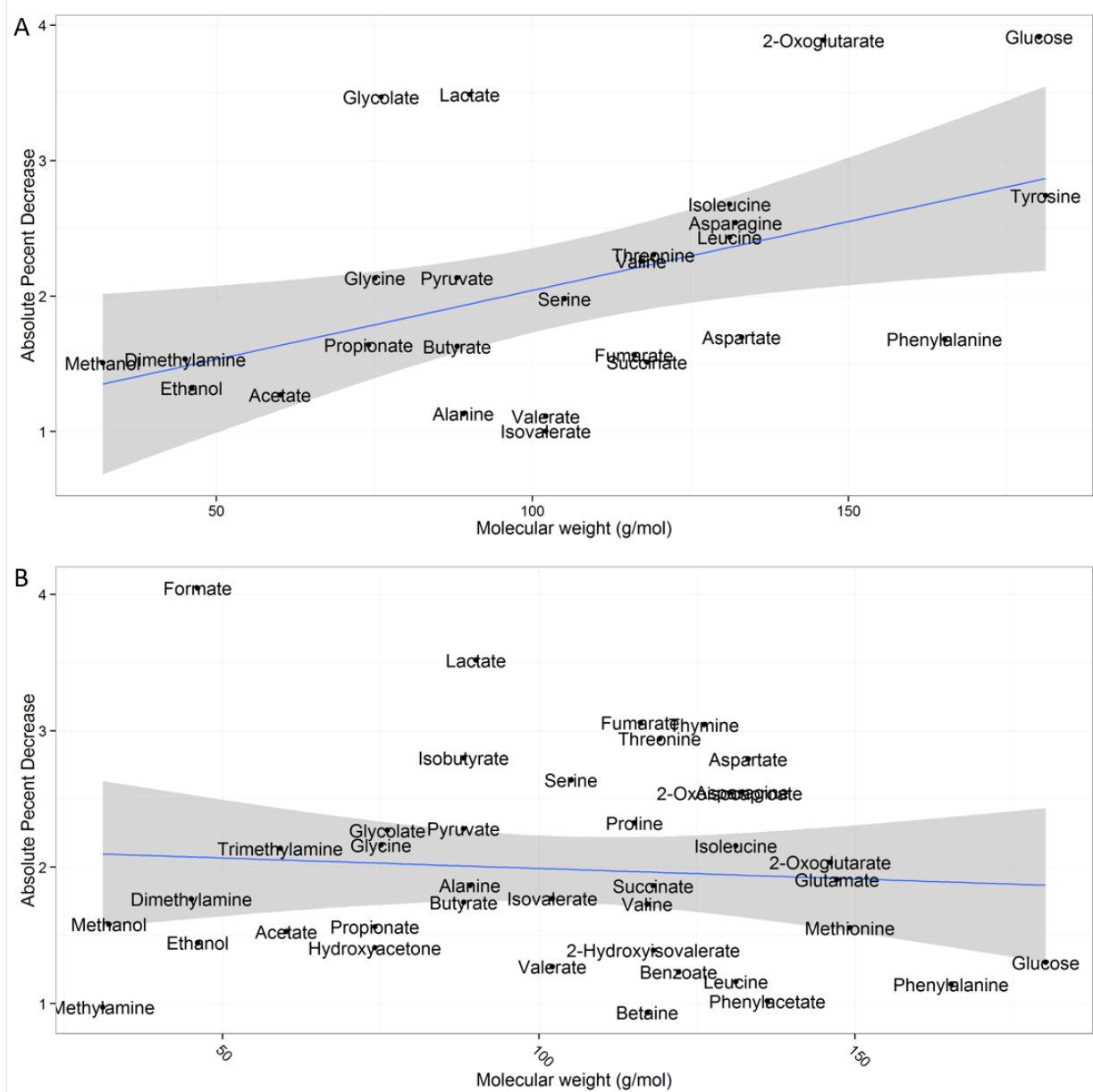


Figure 6-2. Effect of molecular weight on percent concentration decrease. (A) Experiment 1, $y = 21.86x + 59.65$, $p\text{-value} = 0.015$; (B) Experiment 2, $y = -3.84x + 110.76$, $p\text{-value}=0.641$.

6.2.3.5 Comparison of Replicate Experiments

52 compounds were identified and quantified in the metabolite profiles in both experiment replicates. The same compounds were found in the metabolite profiles of both experiment sets. An overview of the two experiments indicated that the two sample sets were very similar to each other, as overlaying the metabolite profiles from the two experiment's samples revealed that the compound concentrations followed similar trends (Figure 6-1). In conducting the experiment twice, the degree to which the ultracentrifuge effects were reproduced from one experiment to the next was assessed. 38 compounds exhibited the same trend in both experiments, in both the concentrations observed and the change in concentrations observed (Table A5). In some compounds however, the concentrations observed were biased by the experiment but the rate of concentration change remained the same between Experiment 1 and Experiment 2 (Figure 6-3). There was no consistency across these 14 compounds which had greater concentrations in the second experiment. Furthermore, five of these compounds, benzoate, beta-alanine, betaine, pyroglutamate and uracil, did not have significant changes in concentration with ultracentrifuge speed in either experiment. Even though the experiment bias observed in these compounds were significant, the importance of this effect was revealed when the magnitude of the experiment difference was examined relative to the concentrations observed (Table 6-2). Most compounds with significant bias only exhibited concentration differences between 2-10% of the observed concentration. Uracil, beta-alanine and pyruvate changed by $11 \pm 2 \%$, $15 \pm 2 \%$ and $19 \pm 1 \%$, respectively, and of these three, only pyruvate demonstrated a significant decrease with ultracentrifugation.

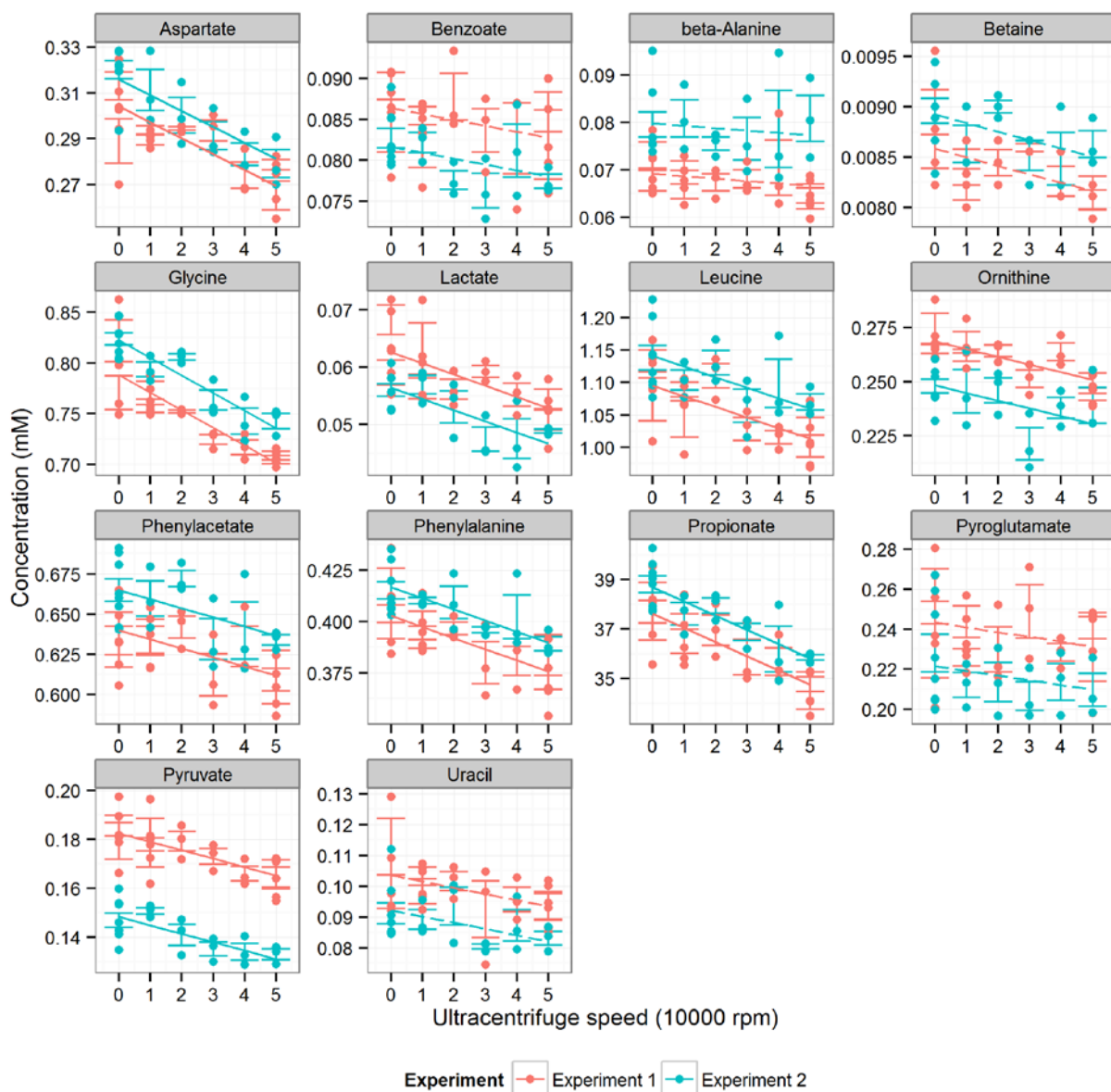


Figure 6-3. Compounds that had significant experiment bias based on MLR model. Slopes significantly different than zero shown as solid lines, otherwise, depicted as dashed lines.

Table 6-2. Absolute percent change due to experiment bias based on: $y = \beta_{speed}x_{speed} + \beta_{experiment}x_{experiment} + intercept$.

Compound	Intercept	Experiment Bias	Bias Standard Error	Absolute Percent Change	Standard Error of Percent Change
Aspartate	0.30	0.012	0.003	3.9	1.1
Benzoate	0.09	-0.005	0.001	5.5	1.4
beta-Alanine	0.07	0.011	0.002	15.4	2.6
Betaine	0.009	0.0003	8.88E-05	4.0	1.0
Glycine	0.79	0.0347	0.005	4.3	0.7
Lactate	0.06	-0.006	0.001	9.0	2.0
Leucine	1.1	0.046	0.012	4.2	1.1
Ornithine	0.27	-0.02	0.003	7.5	1.1
Phenylacetate	0.64	0.025	0.005	3.9	0.9
Phenylalanine	0.40	0.014	0.004	3.5	0.9
Propionate	37.6	1.1	0.27	2.8	0.7
Pyroglutamate	0.24	-0.02	0.006	8.9	2.3
Pyruvate	0.18	-0.03	0.002	18.7	1.2
Uracil	0.10	-0.01	0.002	11.0	2.4

6.2.4 Discussion

6.2.4.1 Profiling Variability

The metabolite profiles obtained in these experiments can be relied upon to provide precise values. The eight samples subjected to single-step filtration in Experiment 2 were used to evaluate the variability of each compound concentration. All of the compounds demonstrated good precision in their quantification, with RSD values less than 15%. The only two compounds that exhibited greater spread than that were p-cresol and thymine, which were near their limits of quantification. These guidelines for acceptable thresholds of precision are outlined in bioanalytics methods validation protocols from the American Food and Drug Administration (FDA)¹³².

6.2.4.2 Pre-Processing Treatments

One of the elements tested involved the presence or absence of a “soft spin” step prior to any ultracentrifugation. The purpose of the “soft spin” was to remove cells and cell debris without exposing the supernatant to intracellular metabolites. This step was excluded for two ultracentrifugation conditions in order to isolate the effects of the “soft spin.” When the metabolite concentrations from duplicate experiments were evaluated as a function of all the experiment factors (ultracentrifuge speed, pre-processing treatment and experiment replicate) and their possible interactions, pre-processing treatment did not influence the concentration observed for any of the compounds. Therefore, the pre-processing step of performing a soft spin prior to ultracentrifugation did not affect the metabolite concentrations.

6.2.4.3 Effect of Ultracentrifuge Speed

Negative linear relationships between compound concentration and centrifugal force applied to the samples were consistently observed throughout the metabolite profiles, prevalent to 37 of the 52

compounds observed. While the relationship between compound concentration and centrifugal force is clear, the absolute rate of concentration change does not seem to be a function of the molecular weight of the compounds. Furthermore, the concentration changes observed were minimal, averaging ~2% decrease in concentration per 10 000 rpm. The only compound found to increase linearly with centrifugal force was p-cresol, averaging at 45.4% per 10 000 rpm. Given the high rate of increase with centrifugal force applied, compounded with the contextual relevance this compound has to enteric microbial metabolism, it is evident that importance should be placed on being transparent regarding sample preparation procedures. Anaerobic catabolism of aromatic amino acids, phenylalanine, tyrosine and tryptophan, occur at the aliphatic side chains to produce phenols like p-cresol. p-Cresol has been shown to be produced by *Clostridium difficile* and *C. scatologenes* from tyrosine via decarboxylation of 4-hydroxyphenylacetate^{133,134}. Furthermore, p-cresol was formed from cell-free extracts of *C. difficile*, where decarboxylation of 4-hydroxyphenylacetate was facilitated by combining the protein fraction of the cell-free extracts with the filtrate fraction, serine, or an α -keto acid like pyruvate¹³⁴. The proposed mechanism for converting p-cresol from 4-hydroxyphenylacetate was a two-protein system where a decarboxylase was kept in the active, reduced form by an α -keto acid dehydrogenase. While 4-hydroxyphenylacetate was not observed in the samples, it has been observed in the metabolite profiles of other feces-derived cultures within this body of work. Therefore, it is possible that the level of 4-hydroxyphenylacetate in this sample set was below the level of detection. Upon examining the absolute rate of concentration change with ultracentrifuge speed, the combined rate of tyrosine and phenylalanine was still ~2 times lower than that of p-cresol, suggesting that if tyrosine and phenylalanine does play a role, they cannot be the only source. Tryptophan was not observed in the metabolite profiles.

6.2.4.4 Comparison of Replicate Experiments

Samples obtained from the same pool of bioreactor effluent were used in both replicate experiments to eliminate sample-to-sample variation. Indeed the compounds that were identified in the metabolite profile from one replicate experiment to the next were identical. The concentration changes observed in the duplicate experiments were explained independently by ultracentrifuge speed and experiment. This is important because an interaction effect between these two factors, indicated by different rates of change in concentration depending on the experiment, would infer that compounds were inconsistent in their response to ultracentrifugation from one experiment to another, even under controlled conditions. The linear regression modelled above demonstrated that there was no significant interaction occurring between any of the experimental factors, so the metabolites' response to the sample preparation protocols was consistent, and thus, predictable. Experiment bias, indicated by a parallel shift in the linear regression from one experiment to another, would mean that there were overall concentration differences between experiment sample sets. Of the 14 compounds that had significant experiment bias, most of these compounds differed by only 2-10% of the observed corresponding compound concentration. Uracil, and beta-alanine had greater discrepancies between experiments, revealing that these two compounds may be sensitive to factors not controlled within this experiment. Pyruvate levels were also not reproduced from one experiment to the next, as it also demonstrated significant experiment bias. However, the impact of ultracentrifuge speed on pyruvate levels was reproduced, with levels decreasing by 2% per 10 000 rpm in both replicate experiments. In general, the compound concentrations between the two sample sets were identical and the compounds response to ultracentrifugation could be predicted.

6.3 Effects of Varying Inoculum Biomass Proportions

6.3.1 Experiment Objectives

The objective of this experiment was to determine the impact of the biomass makeup of inoculants on the resulting *in vitro* community. To evaluate this, defined communities were cultured from inocula that contained the same bacterial species, but the proportional makeup of the bacteria in the seeding culture varied across inocula. The hypothesis driving this experiment was that *in vitro* communities develop similar metabolite profiles when the seeding culture is composed of the same bacterial strains, regardless of the proportions at which the bacteria are present in the inoculum.

6.3.2 Materials and Methods

6.3.2.1 Experiment Outline

Two different defined communities, Microbial Ecosystem Therapeutic-2A (MET-2A) and MET-3A were studied in a CSTR system that consisted of twin vessels, operating in parallel. In a given CSTR run, one vessel was inoculated with a defined community and the other was inoculated with the same defined community, but with different biomass proportions of the same bacteria strains. The varied communities were termed MET-2B or MET-3B, depending on the respective parent community. Culturing each defined community and its reciprocal community was repeated 5 months after the initial experiment (Figure 6-4). Therefore, Run 1 and 3 are duplicate experiments of each other, and Run 2 and 4 are duplicate experiments.

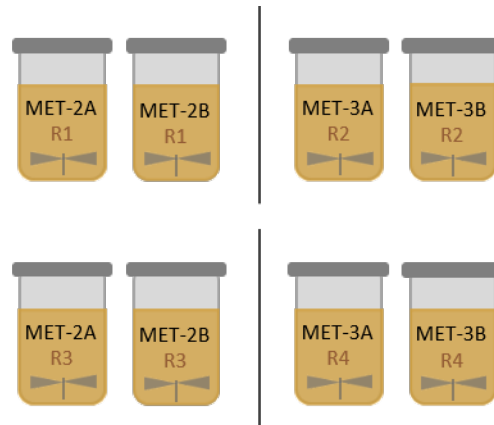


Figure 6-4. Experiment overview. Variants of two different defined communities derived from fecal communities obtained from different donors were examined: MET-2 and MET-3. The A and B variants of the MET-2 and MET-3 communities have the same bacterial strains, but have varied proportional makeup of those bacteria. Run 1 (R1) and Run 3 (R3) were duplicate experiments and Run 2 (R2) and Run 4 (R4) are duplicate experiments.

6.3.2.2 Sampling and Storage

Samples were collected directly from the bioreactor vessel throughout the duration of the bioreactor runs to monitor the cultures' growth. Long, sterile sampling needles that punctured a silicone septum were used to directly access the bacterial culture in the bioreactor vessel. 4-mL samples were aliquoted into 2 2-mL tubes and stored at -80 °C. Samples from day 6, 14, and 20 from each bioreactor run were selected for metabolite analysis. Samples were thawed, then centrifuged at 14 500 rpm for 5 min to pellet the cell fraction. The supernatant was subsequently filtered through 0.45 µm and 0.22 µm syringe filters (Corning; Corning, New York) and stored for up to 4 days at 4 °C until NMR scanning.

6.3.2.3 Statistical Analysis

The metabolite profiles of each community culture were annotated by identifying the compounds' chemical taxonomic classes. The metabolite profiles could then be analyzed according to these taxonomic classes, examining the prevalence of each compound class within a profile, and how these class proportions change when weighted to compound concentrations. The chemical taxonomic class of each compound was identified based on classifications listed in the Human Metabolome Database¹³⁵ (www.hmdb.com). The magnitude of the metabolite profile represented by a given compound class, referred to as the class weight, was calculated as the sum of concentrations of metabolites of that class within a metabolite profile. The class weight as a proportion, referred to as class percentage, was presented by dividing the class weight by the total sum of metabolite concentrations in the profile. The mean class weight within a given culture was also calculated. Two samples were selected at random to be scanned in triplicate and be used to assess profiling variability. The precision of the compound concentrations were assessed based on the relative standard deviation (RSD) in each triplicate set. RSD was calculated as shown in Appendix B1. The mean of each set of

triplicates was used as the concentration for each respective sample (that were analyzed in triplicate). Student's t tests were also applied to assess the run-to-run differences between the two pairs of duplicate experiments performed for each defined culture, MET-2A, MET-2B, MET-3A and MET-3B. It should be noted that the three time-course data points from each run were considered as one group when comparing the duplicate runs. While a time-matched comparison across duplicate runs would be more ideal, comparing the runs across all three time points was sufficient as the two runs were expected to behave the same. Even so, a more stringent level of significance would give more allowances for some dynamicity in the progression of the runs. Moreover, performing the Student's t test repeatedly (for every compound and every defined culture) leaves the test to be prone to type I error. For these two reasons, a Bonferroni correction was applied to the alpha value, which defined the significance level as $p \leq 0.0075$ (Appendix B3.). PCA was used to analyze the metabolite profiles. This multivariate analysis was performed using the function "prcomp" within the R base package "stats". Linear regression was used to assess the trends of compound concentrations over time (day 4-20) via the "aov" function in the R base package "stats". This function is a form of ANOVA that assesses the variation of concentrations about a linear regression model and compared it to the variation of the concentrations about the grand mean. Significance was determined at the 90% confidence level ($p \leq 0.1$). This confidence level was chosen because in this case, several weak linear trends had significance values between 0.05 and 0.1.

6.3.3 Results

6.3.3.1 Profiling Variability

For the majority of compounds, the relative standard deviation was less than or equal to 15%. (Figure A5). Compounds that had a spread greater than 15% from the mean are shown in Table 6-3.

However, even with the spreads observed, the two samples done in triplicate could still be distinguished from each other as demonstrated by PCA (Figure 6-5).

Table 6-3. Compounds with RSD in replicate samples > 15%

Sample 1		Sample 2	
Compound	RSD (%)	Compound	RSD (%)
Betaine	16.5	Benzoate	19.3
Fructose	24.3	Fructose	21.4
Fumarate	21.3	Glucose	21.4
Hydroxyacetone	16.8	Glycine	28.0
Proline	28.4		

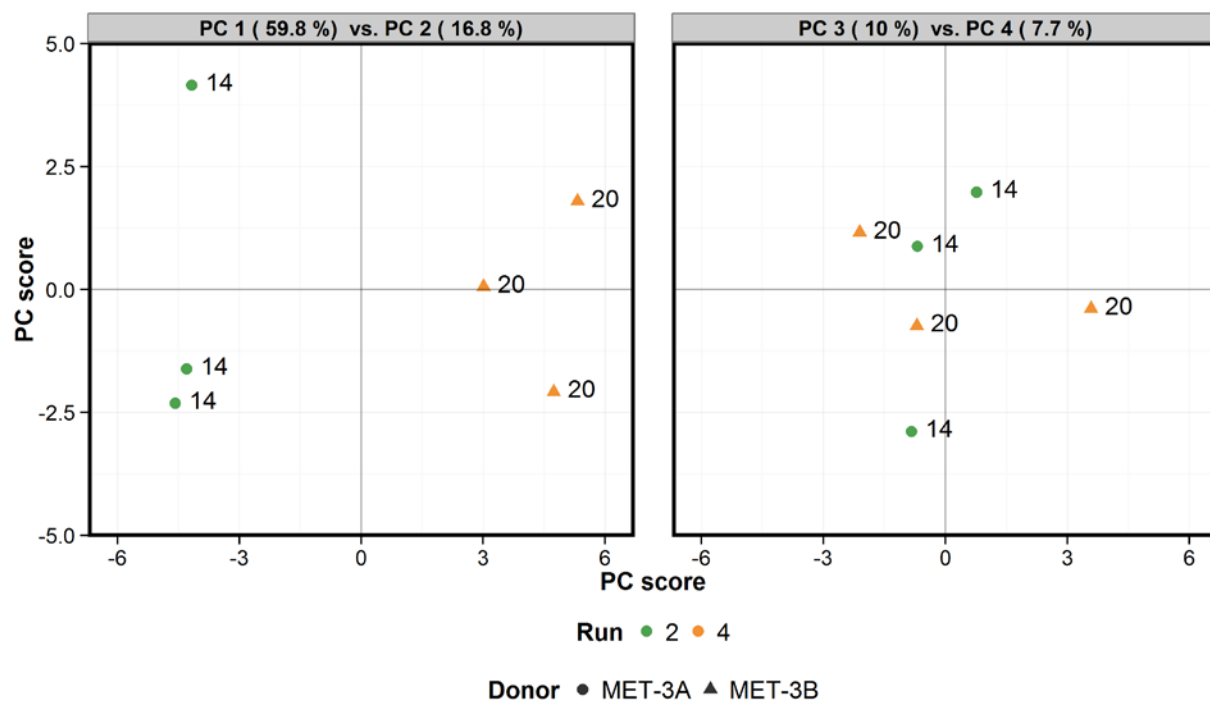


Figure 6-5. PCA score plot of triplicate samples.

6.3.3.2 Replicate Runs

Run 1 and Run 3 were duplicate experiments that propagated MET-2A and MET-2B. Run 2 and Run 4 were duplicate experiments that propagated MET-3A and MET-3B. While the duplicate Runs may have been identical in the seeding population, there were slight differences in the metabolite profiles of each pair of replicates (Figure A6). These differences became apparent in a PCA score plot of all the samples (Figure 6-6). In this overall model, Run 2 and 4 produced similar metabolite profiles of one another, but Run 1 and 3 clustered distinctly from one another. The variation along principal component (PC) 1, along the x-axis, accounted for 60.3% of the total variation and accounted for the differences between the defined culture types, MET-2A and MET-3A. The variation along the y-axis, PC2, occurred predominantly due to differences between Run 1 and Run 3 and differences between Run 2 and 4. Since differences between the duplicate runs occurred along PC2, which accounted for only 8.3% of the total variation, the differences between the duplicate runs were much smaller than the differences between the types of defined communities (MET-2A vs MET-3A).

Since all four cultures were propagated in duplicate runs, the reproducibility of each culture could be assessed by using Student's t-test to compare each compound between the two runs. The compounds that differed significantly between the duplicate runs are shown in Table 6-4. All of these compounds, except for fumarate, had within-run mean concentrations greater than 15% of the between-run mean concentration.

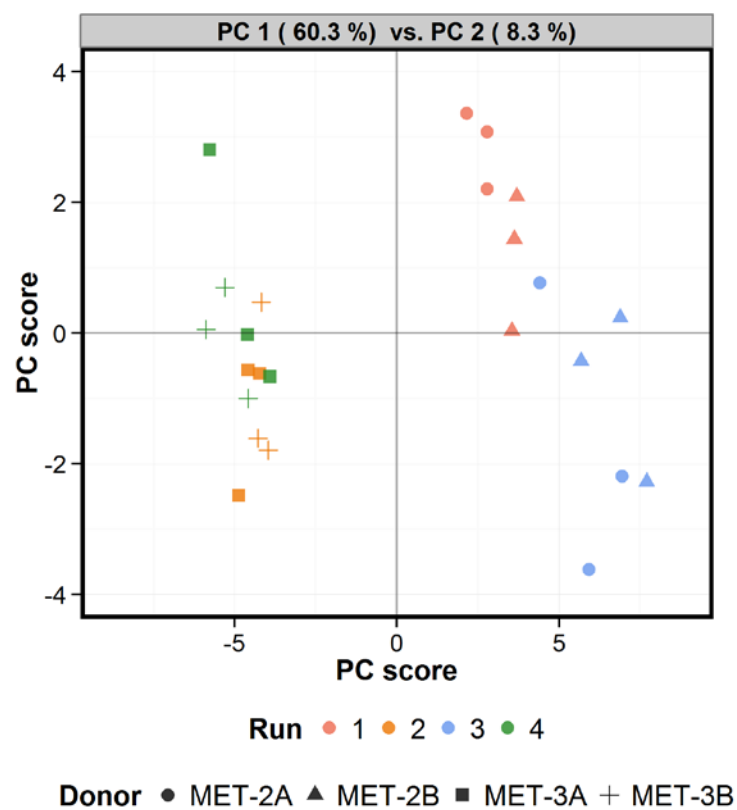


Figure 6-6. PCA score plot of metabolite profiles of all bioreactor samples. Metabolite concentrations have been mean centered and unit variance scaled in order for equal-weighted comparison between compounds.

Table 6-4. Compounds that differed significantly between duplicate runs.

Compound	Culture	T Statistic	P value	Absolute Percent Difference ³
Fumarate	MET-2A	-5.1	0.008	8.8
Propionate	MET-2A	5.3	0.007	28.2
Aspartate	MET-2B	-10.7	0.002	28.5
Glycine	MET-3A	5.2	0.007	32.3
Glycine	MET-3B	5.5	0.005	43.9
Isoleucine	MET-3B	13.2	0.0003	16.5

³Percent difference = |mean concentration within a run – mean concentration from duplicate runs|/mean concentration from duplicate runs.

6.3.3.3 Metabolite Profiles

Metabolite profiles were constructed for 24 samples. Each profile contained 41 compounds where the assignment of the signal's compound identity and concentration were considered to be high confidence. The total concentration of the metabolite profiles was greater in MET-3A-based cultures (MET-3A and MET-3B) than to MET-2A-based cultures (MET-2A and MET-2B), and similar between a defined culture and its reciprocal (Figure 6-7). The structure of the metabolite profiles was fairly consistent across all four different seeding populations. Based on the number of compounds profiled, amino acid and amino acid derivatives (39%) comprised most of compounds identified, followed by fatty acids and conjugates, carboxylic acids and derivatives and monosaccharides which each accounted for ~10%. When the taxonomic classes in the metabolite profiles were weighted to the compounds' concentration, the proportions of classes were drastically altered. Carboxylic acids and derivatives dominated in the profiles, followed by fatty acids and conjugates, alcohols and polyols and amino acids and derivatives (Table 6-5). The distribution of these weighted compound classes was fairly similar across all four cultures, with slight differences between the MET-2A-type cultures and the MET-3A-type cultures. The most noticeable differences were a dramatic increase in carboxylic acids and a decrease in fatty acids in MET-3A-type cultures compared to MET-2A-type cultures.

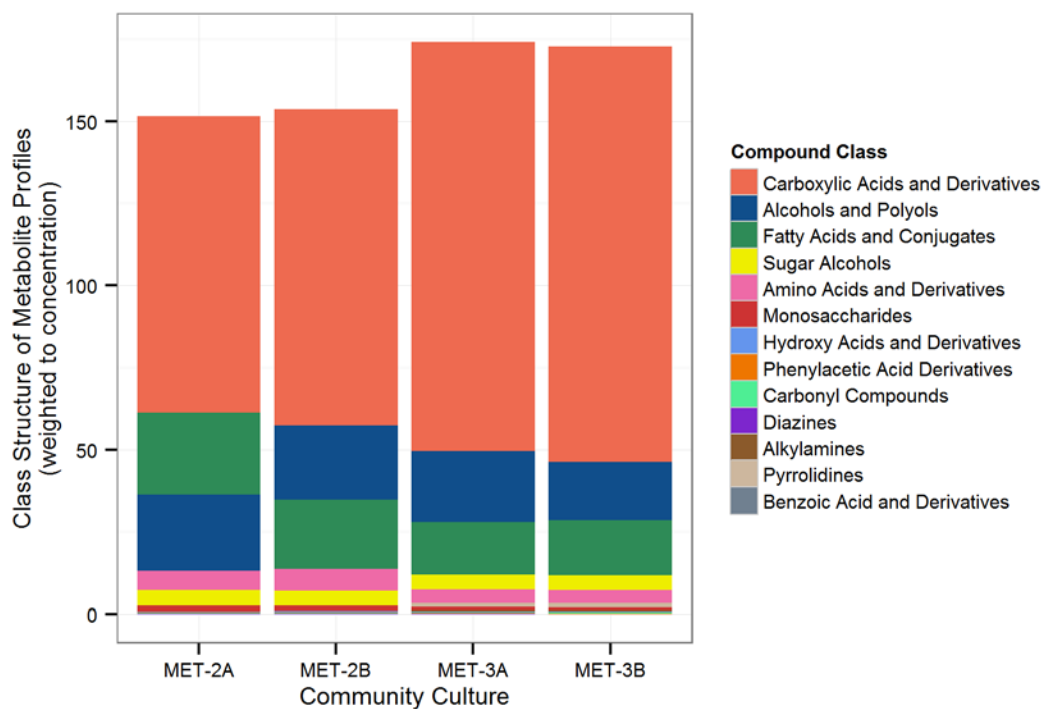


Figure 6-7. The breakdown of compound classes in metabolite profiles of each defined culture. Proportions of compound classes were weighted to compound concentration. The class weights shown are the mean class weights of profiles within a given culture.

Table 6-5. Proportion of Compound Classes, weighted to concentrations.

Compound Class	MET-2A	MET-2B	MET-3A	MET-3B
Carboxylic acids and derivatives	59.4	62.5	71.4	73.0
Fatty acids and conjugates	16.5	13.6	9.2	9.9
Alcohols and polyols	15.3	14.8	12.4	10.3
Amino acids and derivatives	3.9	4.4	2.3	2.4
Sugar alcohols	3.1	2.9	2.6	2.5
Monosaccharides	1.1	1.1	0.6	0.5
Hydroxy acids and derivatives	0.2	0.2	0.2	0.2
Phenylacetic acid derivatives	0.1	0.2	0.1	0.1
Carbonyl compounds	0.1	0.01	0.1	0.1
Diazines	0.09	0.01	0.1	0.1
Alkylamines	0.07	0.08	0.2	0.2
Benzoic acid and derivatives	0	0	0.01	0.01
Pyrrolidines	0	0	0.7	0.7

6.3.3.4 Defined Cultures vs Their Reciprocals

Most of the compounds tend to a common point over time between a given culture and its reciprocal. Compounds that reach this common point at Day 20 in a weak linear fashion in at least one run are shown in Figure 6-8A and Figure 6-10A. The compounds shown in Figure 6-8B and Figure 6-10B reach a common point at Day 20, but in a non-linear pattern. It is notable that the achievement of a common point is not always observed across duplicate runs (Figure 6-9 and Figure 6-11). The compounds that do not reach a common point by Day 20 in any of the runs are listed in Table 6-6.

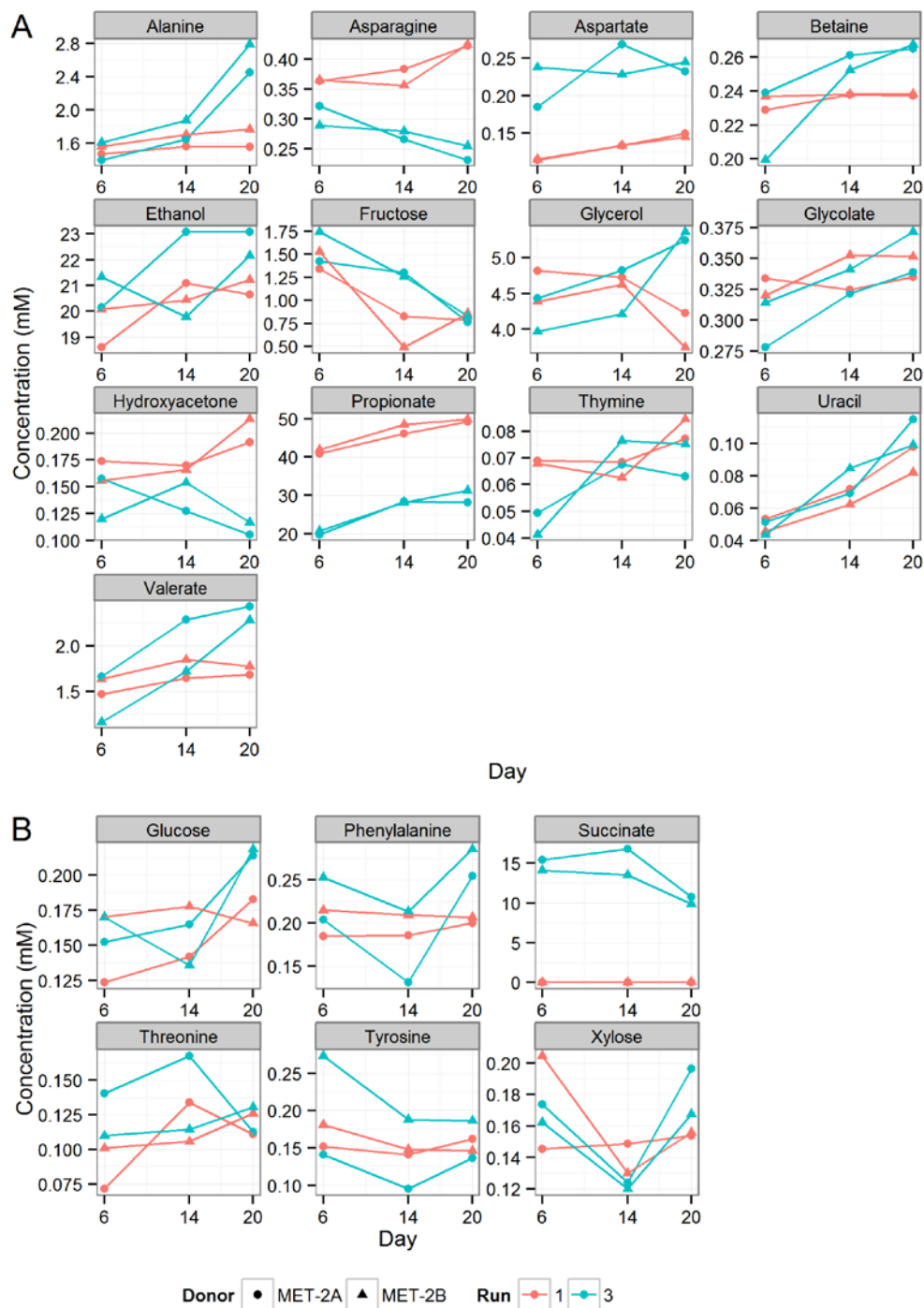


Figure 6-8. Compounds that reach the same level by Day 20 across a MET-2A culture and its reciprocal (MET-2B). (A) Compounds that display time trends that demonstrated a linear model in both runs. (B) Compounds that have time trends that are not linear.

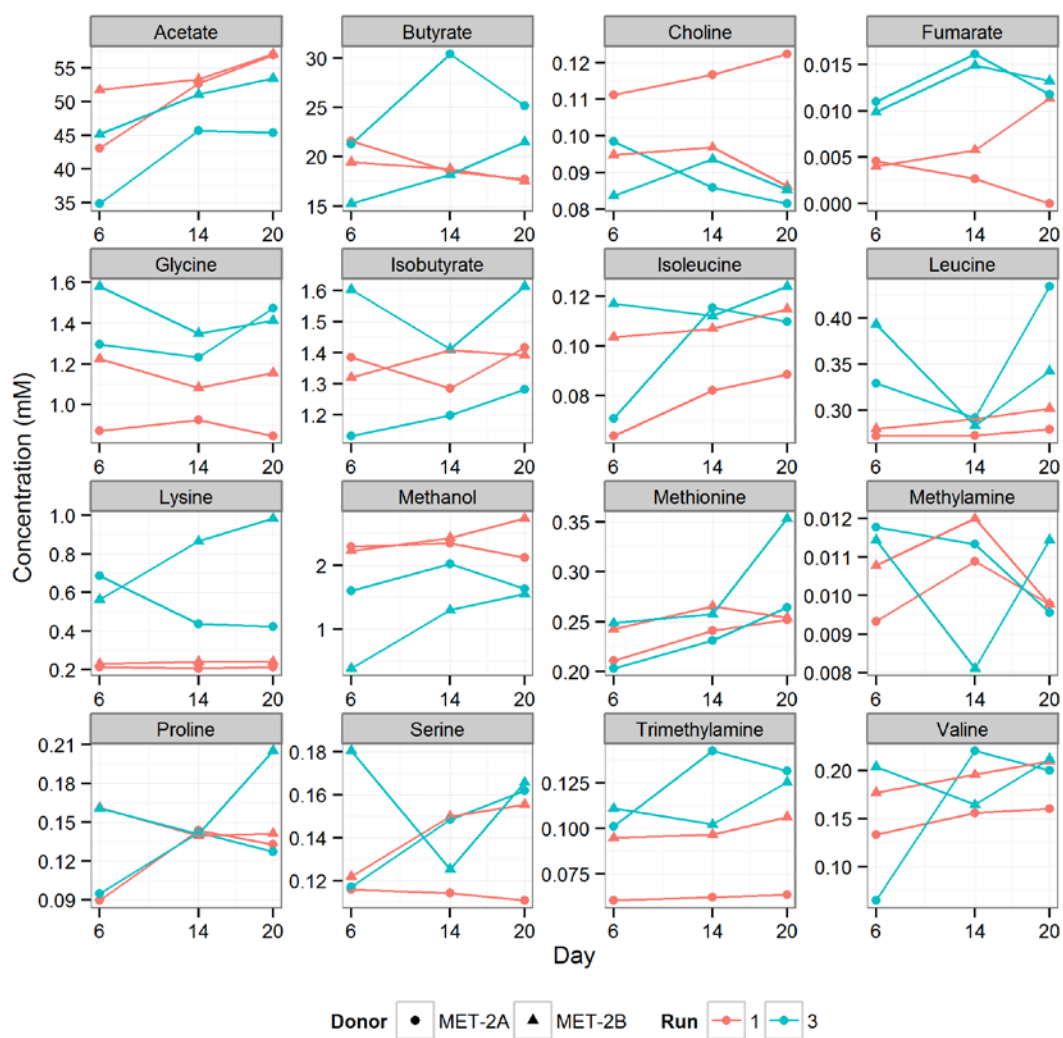


Figure 6-9. Common point reached in only one run for MET-2A and MET-2B.

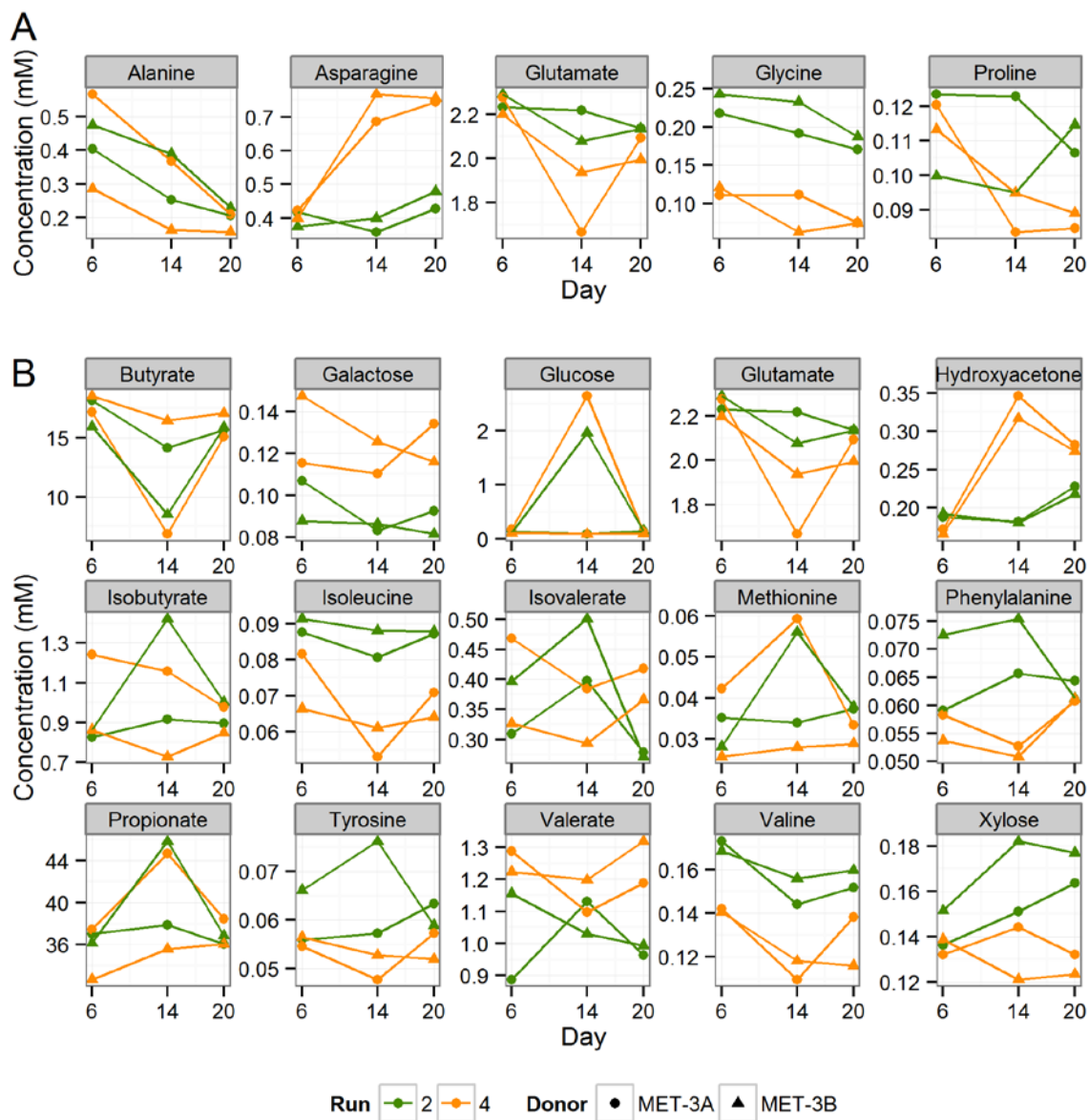


Figure 6-10. Compounds that reach the same level by Day 20 across a MET-3A culture and its reciprocal (MET-3B). (A) Compounds that display time trends that demonstrate a weak linear model in at least one of the runs. (B) Compounds that have time trends that are not linear.

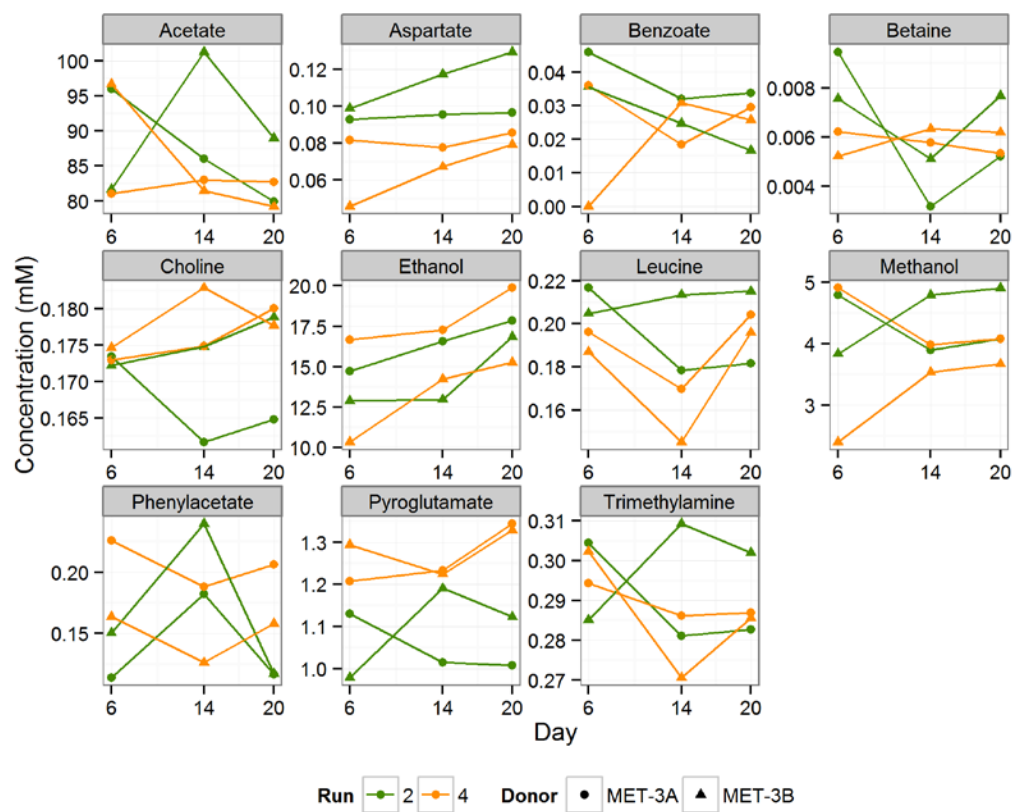


Figure 6-11. Common point reached in only one run for MET-3A and MET-3B

Table 6-6. Compounds that do not reach a common point between a culture and its reciprocal by Day 20 in either of the duplicate runs.

MET-2A vs MET-2B	MET-3A vs MET-3B
Galactose	Fructose
Isovalerate	Fumarate
Phenylacetate	Glycerol
	Glycolate
	Methanol
	Methylamine
	Serine
	Thymine

6.3.4 Discussion

6.3.4.1 Profiling Variability

Assessment of the metabolite profiles of replicate samples revealed that most compounds were quantified with a relative standard deviation of less than or equal to 15%. Since FDA regulations on bioanalytical method validation recommend 15% RSD, or 20% if at the lower limit of quantification, as the tolerance limit in assessments of precision¹³², the metabolite profiles can be considered to have relatively little variation. The compounds that demonstrated more spread tended to be compounds with only a few number of clusters in its NMR signal, a higher degree of convolution with other spectral signals, or the compound peaks were masked due to poor signal-to-noise ratios. Even though these, or any combination of these spectral characteristics can increase the variation in quantification, the precision was still high enough to be easily distinguished from other samples.

6.3.4.2 Replicate Runs

The run-to-run effects comment on the repeatability of culturing the same bacterial population. The multivariate model that described the metabolite profiles explained that only 8.3% of the total variation could be attributed to run-to-run variability. Comparing the metabolite profiles of each of the duplicate runs revealed that most compounds levels were reproduced across runs, and only one or two compounds had significantly different concentrations from one replicate run to the next.

It should be noted that all time points from the growth of a culture were included in the run-to-run comparisons since only one data point was available for each time point. A more accurate evaluation of repeatability would compare runs only at time-matched points, where each time point had at least 3 data points. The implications of grouping all three time-points together in making run-run comparisons was that the test is more sensitive to compounds that remain constant over time

within one run. Compounds that fluctuate largely between day 4 and day 20 of a given CSTR culture would be more likely identified as an insignificant difference with other culture runs. That being said, the replicate runs can still be differentiated by gross concentration differences in the metabolite profiles.

6.3.4.3 Metabolite Profiles

The metabolite profiles of the four bacterial communities were very similar to one another, with the greatest degree of similarity between a culture and its respective reciprocal. The same 41 compounds were observed in the metabolite profiles in all four cultures. The most apparent differences in metabolite profiles were evident when comparing the profiles between MET-2A-based cultures and MET-3A-based cultures. Overall, the total concentration of compounds profiled was greater in MET-3-based cultures (MET-3A and MET-3B) compared to MET-2-based cultures (MET-2A and MET-2B). Furthermore, there was a larger proportion of carboxylic acids in the profiles of MET-3-based cultures (~72%) than in the MET-2-based cultures (~60%). This difference was mostly attributed to acetate, where the average concentration observed in MET-3A-based cultures was 86 ± 2 mM standard error, while the average concentration observed in MET-2A-based cultures was 49 ± 3 mM standard error. The comparatively low concentrations of fatty acids and conjugate compounds in MET-3-based cultures (~9% in MET-3-based cultures and ~15% in MET-2-based cultures) was mostly due to lower levels of butyrate. The mean concentration of butyrate in MET-3A-based cultures was 15 ± 1 mM standard error compared to 21 ± 2 mM standard error found in MET-2A-based cultures. Concentrations of valerate and isovalerate were also considerably lower in MET-3A-based cultures but they were found at fairly low concentrations in all four cultures, so they only accounted for ~1-2% of the metabolite profiles.

6.3.4.4 Defined Cultures vs Their Reciprocals

In general, there was very little distinction between each defined culture and their respective reciprocal culture. In fact, an overview of the time course progression of the compound changes indicated the cultures and their biomass-varied counterparts reached similar metabolite profiles by the 20th day post-inoculation. For some compounds, this progression happened linearly, while in others, the compounds experienced some fluctuation before reaching a common concentration point at day 20. Only a select few compounds were not observed to reach a common point between a defined culture and its reciprocal cultures in either duplicate runs. While the convergence of compound concentrations by day 20 between a defined culture and its reciprocal culture was a first indication that the defined cultures tended towards a single metabolic steady state regardless of the biomass of species isolates in the inoculum, there are still aspects of this experiment that must be further investigated. First of all, it should be noted that the linearity was evaluated based on 3 time points of a given culture. For more robust evaluations on time series patterns, more time points should be evaluated, with multiple samples of each time point. This would enable more accurate modelling of the cultures' metabolic activities. Moreover, previous investigations reported that fecal cultures grown in CSTR systems do not stabilize in terms of species composition until ~5 weeks after inoculation⁵. Therefore, the cultures should be observed for longer in the each bioreactor run to observe when a true metabolic steady state is actually reached, and to determine if the similarities in the metabolite profiles of a culture and its reciprocal would maintained.

Chapter 7 Metabolomic Analysis of Human Fecal Microbiota Propagated *In Vitro* in a CSTR System

7.1 Chapter Objective

The purpose of this chapter was to demonstrate ways in which the metabolomic techniques and analyses used thus far can be applied to gain further understanding of the human gut microbiota. In the first experiment, metabolite profiles were used to describe feces-derived and defined *in vitro* communities. This study exemplified the potential of performing metabolomic analyses on *in vitro* gut microbial communities, as unique metabolic functions were identified in each bacterial culture. The second experiment supported the application of this metabolomic analysis for studying the human gut microbiota in a clinical and pharmaceutical framework. Several *in vitro* gut microbial communities were subjected to various types of perturbations, including antibiotics, a hormone and a defined community culture used to challenge the bioreactor culture. Examining the metabolomic responses to a wide range of perturbations set the preliminary groundwork in for raising attention to the metabolomic implications relevant to industrial and clinical research.

7.2 Metabolomic Analysis of the Human Gut Microbiota: Comparison of Feces-Derived Communities and Defined Mixed Communities

7.2.1 Journal Article Authorship

Sandi Yen¹, Julie A. K. McDonald², Kathleen Schroeter², Stanislav Sokolenko¹, Eric J. M. Blondeel¹, Emma Allen-Vercoe^{2#}, Marc G. Aucoin^{1#}

¹Waterloo Institute for Nanotechnology, Department of Chemical Engineering, University of Waterloo,

Waterloo, Ontario, Canada, N2L 3G1

²Department of Molecular and Cellular Biology, University of Guelph, Guelph, Ontario, Canada, N1G 2W1

#corresponding authors: Marc G. Aucoin maucoin@uwaterloo.ca, Emma Allen-Vercoe eav@uoguelph.ca

7.2.2 Justification of Original Work

With minor editorial changes to fulfill formatting requirements, this chapter is substantially as it has been submitted to the Journal of Proteome Research on October 31, 2014. The study was performed in collaboration with the Allen-Vercoe group at the University of Guelph. The partnership was formed where the Allen-Vercoe group brought forth expertise in culturing microbial communities representative of the human gut in an *in vitro* CSTR system and the Aucoin group extracted, analyzed and interpreted metabolomic data obtained from the microbial communities. Preparation and operation of the gut-mimicking CSTR system was performed by Julie McDonald and Kathleen Schroeter. They also collected raw fecal samples, prepared all community inocula, collected

bioreactor and effluent samples, and processed the samples in preparation for metabolomic analysis. 1D-¹H NMR spectroscopy was performed by Sandi Yen. The manuscript was written by Sandi Yen. This work is included in this thesis because it applied the skills and analytical techniques developed from earlier studies presented in this thesis.

7.2.3 Abstract

The extensive impact of the human gut microbiota on its human host calls for a need to understand the types of communication that occur amongst the bacteria and its host. A metabolomics approach can provide a snapshot of the microbe-microbe interactions occurring, as well as variations in the microbes from different hosts. In this study, metabolite profiles from an anaerobic continuous stirred-tank reactors (CSTR) system supporting the growth of several consortia of bacteria representative of the human gut were established and compared. Cell-free supernatant samples were analyzed by 1D-¹H Nuclear Magnetic Resonance (NMR) spectroscopy, producing spectra representative of the metabolic activity of a particular community at a given time. Using targeted profiling, specific metabolites were identified and quantified based on NMR analyses. Metabolite profiles discriminated each bacterial community examined, demonstrating that there are significant differences in the microbiota metabolome between each cultured community. We also found unique compounds that were identifying features of individual bacterial consortia. These findings are important because they demonstrate that metabolite profiles of gut microbial ecosystems can be constructed by targeted profiling of NMR spectra. Moreover, examination of these profiles sheds light on the type of microbes present in the gut and their metabolic interactions.

7.2.4 Introduction

In the wake of metagenomic research aiming to fully describe the gut microbiome, metabolomic research has emerged as a functional counterpart¹³⁶ that provides insight into the metabolic activities engaged by whole communities of microorganisms. Metabolomics is the study of the entire set of compounds, or metabolites, produced by a given organism, referred to as a metabolome. As a reflection of enzymatic activities and metabolic pathways carried out by the organism¹³⁷, the application of metabolomics to whole community bacterial cultures confers the ability to take a “snapshot” of the complex microbial population. Within the context of the human gut microbiota, the mixed culture metabolome not only reflects the metabolic activities of the enteric bacteria, but it also represents the bacterial compounds to which the human body is exposed.

Metabolomic studies of the gut microbiota often use metabolites from fecal extracts to compile metabolomes that capture the gut bacteria-host interaction^{2,63,138}. The metabolomic data obtained directly from stool samples contains metabolites from both host and microbial activities. Thus, while fecal metabolomics provide insight into the gut microbial-host co-metabolism, it is not the ideal source for studying the independent metabolic contributions from either the host or enteric bacteria. Furthermore, the composition of the gut microbiota varies depending on host age, nutrition, behaviour and health^{52,53,139}. Therefore, while fecal metabolomic studies bring forth understanding of the microbial-mammalian metabolic axis, the exact conditions are difficult to recreate to further investigate potential relationships. Taking the approach of microbial metabolomics to study the gut microbiota, analysis of the extracellular metabolites (the ‘exometabolome’) produced by the enteric bacteria circumvents host biology variability and focuses the metabolomic analysis on the community’s metabolic state.

The robustness of using metabolomic approaches to study the gut microbiota relies heavily on the acquisition of high-quality data. The most common analytical platforms used in quantifying metabolomes are mass spectrometry (MS) and nuclear magnetic resonance (NMR) spectroscopy¹⁴⁰. Both platforms have advantages and limitations that influence for which situations the analytical techniques are best suited. The methodology and application of MS and NMR in metabolomics have been extensively reviewed^{23,140–142}. Within the scope of compiling enteric microbial metabolite profiles, ¹H NMR spectroscopy is the platform of choice because it requires minimal pre-processing of samples, enabling rapid yet reproducible spectroscopic results^{29,39,143}. Its signals provide structural information about the compounds, facilitating identification of known compounds and classification of unknown compounds. Furthermore, targeted profiling techniques have been developed for deconvolution of NMR spectra and provide a level of quantification and identification that was not previously possible with binning algorithms, the traditional tool in analyzing NMR spectra^{31,41}.

In this study, we employed metabolomic characterization of feces-derived bacteria propagated as various mixed cultures in a continuous stirred tank reactor (CSTR) that mimicked the environment of the distal regions of the human gastrointestinal tract^{5,144}. With this approach, we built quantitative metabolite profiles of different gut bacterial communities. The purpose of this study was to describe the metabolome of bacterial communities propagated in a CSTR, and to evaluate identifying characteristics of different populations. In deciphering the ecological dynamics of the gut microbiota, further understanding of the interaction between the gut microbial ecosystem and the human body can be achieved.

7.2.5 Materials and Methods

7.2.5.1 Experiment overview

The Research Ethics Board of the University of Guelph approved this study (REB#09AP011). This study was performed over 5 different bioreactor runs, where a run had one or more bioreactor vessels operating (Figure 7-1). All vessels within a run were inoculated from the same seeding sample. Conversely, each bioreactor run used different seeding sources, whether it was from a different donor, or from the same donor. Three healthy donors provided fresh fecal samples: donor A (male, 44 years-old), donor B (female, 26 years-old) and donor C (male, 25 years-old). None of these donors had a recent history of antibiotic treatment or gastrointestinal disorders within the 9 months prior to collection of their samples. Runs 1 and 2 were inoculated with fecal extracts from donor A whereas Runs 3 and 4 were inoculated with fecal extracts from donor B and C respectively. Run 5 was inoculated with a defined mixed culture, called Microbial Ecosystem Therapeutic-2 (MET-2), which was a defined community consisting of strains derived from donor A's feces (Table 7-1). Runs 1, 2, 3, and 4 were maintained for 48, 53, 41 and 53 days respectively, while Run 5 was maintained for 27 days. Based on previous studies using this system the community cultures reached steady-state as early as 4 weeks after inoculation⁵.

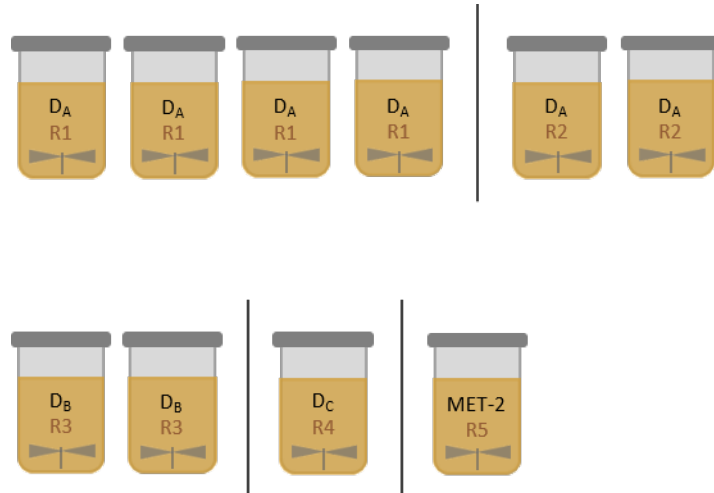


Figure 7-1. Experiment overview. CSTR Bioreactor vessels were inoculated with a seeding culture that were derived from either a human fecal samples, or a defined mixed culture. Fecal samples were obtained from donors A, B and C, as denoted by D_A , D_B , and D_C respectively. A cocktail of defined seeding cultures were used, termed Microbial Ecosystem Therapeutic-2 (MET-2). Bioreactors that were operated in parallel are considered to be part of a single Run (denoted as R1, R2 etc.), where the seeding culture in each bioreactor within a Run are derived from the same inoculant. Samples were collected from bioreactor effluent at various time points during bioreactor operation.

Table 7-1. A list of cultured bacterial isolates from the feces of donor A included in the MET-2 defined community.

Higher taxonomic group	Closest species match, inferred by alignment of 16S rRNA sequence to GreenGenes database	% identity to closest match	Relative abundance ⁴ (by biomass)
Actinobacteria	<i>Bifidobacterium longum</i>	100	1
	<i>Collinsella aerofaciens</i>	100	0.5
	<i>Microbacterium schleiferi</i>	99.34	1.5
	<i>Aldercreutzia equolifaciens</i>	99.76	1
	<i>Micrococcus luteus</i>	97.04	1
Bacteroidetes	<i>Bacteroides ovatus</i>	100	0.5
	<i>Parabacteroides merdae</i>	100	1
Firmicutes	<i>Eubacterium rectale</i>	100	1
	<i>Blautia luti</i>	98.91	1
	<i>Roseburia hominis</i>	99.04	1
	<i>Roseburia lactaris</i>	95.07	0.5
	<i>Ruminococcus albus</i>	96.96	1
	<i>Eubacterium eligens</i>	96.78	1
	<i>Ruminococcus torques</i> (two different strains)	99.27 100	1 1
	<i>Roseburia intestinalis</i>	100	1
	<i>Eubacterium fissicatena</i>	97.67	1
	<i>Eubacterium ventriosum</i>	97.37	1
	<i>Blautia coccoides</i>	99.85	0.25
	<i>Blautia hydrogenotrophica</i>	100	0.1
	<i>Dorea longicatena</i>	100	0.25
	<i>Dorea formicigenerans</i>	99.49	0.25
	<i>Clostridium ramosum</i>	96.14	0.25
	<i>Eubacterium limosum</i>	99.25	0.75
	<i>Streptococcus thermophilus</i>	100	0.25
	<i>Bacillus simplex</i>	98.7	0.25
	<i>Coprococcus catus</i>	99.19	0.25
	<i>Flavonifractor plautii</i>	96.21	0.25
	<i>Streptococcus mitis</i>	100	0.3
	<i>Phascolarctobacterium</i> sp.	99.85	0.1
Proteobacteria	<i>Escherichia coli</i>	100	0.25
	<i>Parasutterella excrementihominis</i>	100	0.5

⁴Relative abundance nits are presented as the equivalent of 1 X 10⁶ µL loopful of bacterial culture scraped from petri dish plates.

7.2.5.2 Bioreactor Operation and Inoculation

Up to four 500-mL Multifors bioreactors (3 Multifors systems, Infors AG, Switzerland), with working volumes of 400 mL were operated as CSTRs. Conditions were set to mimic the distal human gut (37°C, pH 7, gentle agitation, oxygen-free, and fed with a constant supply of mucin and insoluble starch substrates at a flow rate of ~400 mL/day). Vessels were maintained anaerobic by bubbling nitrogen. Operation and maintenance of the bioreactors were conducted as previously specified⁵. The bioreactors were inoculated with either the defined mixed culture, MET-2, or undefined mixed culture obtained from homogenized fecal samples from human donors. 100 mL of inoculum was added to 300 mL of media to comprise a 10% w/v culture. Cultures were gently agitated pH adjusted to 6.9-7.0. The media was prepared in-house, as previously described⁵.

7.2.5.3 Preparation of Fecal Inocula and Defined Mixed Culture

Fecal inocula were prepared as previously described⁵, where donors provided fresh fecal samples. The MET-2 consortium was developed by the Allen-Vercor group of the University of Guelph as a simplified version of the gut microbiota of donor A, and was composed of 33 different strains of bacteria that were isolated from a fecal sample obtained from donor A (Table 7-1). These strains were isolated through plate culture of diluted fecal samples on various media types (deMann Rogosa & Sharpe agar, Brain Heart Infusion agar (Oxoid), Fastidious Anaerobe Agar (Lab90)), with incubation in an anaerobe chamber (Ruskinn) under an atmosphere of N₂:CO₂:H₂ (90:5:5). Isolated colonies were streak-purified and characterized by Sanger sequencing of their 16S rRNA genes with comparison to the Greengenes database (<http://greengenes.lbl.gov/cgi-bin/nph-index.cgi>). Each inoculum was prepared by scraping biomass from select plate cultures and suspending it in pre-reduced CSTR growth medium (75 mL).

7.2.5.4 Sample Collection and Processing

Bioreactor effluent from Run 1 was pooled over a 10-day period and samples were taken from this pooled material. Samples from CSTR Runs 2-5 were collected on a daily basis; sampling volumes did not exceed 2.5% of the total working volume (so not to impact the chemostatic conditions of the bioreactor). Samples were stored at -80 °C, and thawed prior to analysis in a step wise manner (5 hours in -40 °C, overnight in -20 °C and then on ice). Samples were centrifuged (3800 x g, 10 min) to remove cells, then subjected to ultracentrifugation (14650 x g, 4 °C) to remove media particulates. Finally, samples were filtered using 0.22 µm syringe filters (Acrodisc® Supor membrane, Pall Corporation, UK).

7.2.5.5 NMR Spectra Acquisition and Processing

Samples were diluted to contain 10% (v/v) internal standard in preparation for NMR spectroscopy. The internal standard used was a chemical shift indicator (CSI), a solution of 99.9% D₂O with 5 mM 4,4-dimethyl-4-silapentane-1-sulfonic acid (DSS) and 0.2% (w/v) sodium azide to inhibit bacterial contamination. Sample solutions were transferred to 5-mm glass NMR tubes (NE-UL5-7, New Era Enterprises Inc., Vineland, NJ, USA) and inserted into a Bruker Avance 600.13 MHz spectrometer with a Triple Resonance Probe (TXI 600) for scanning. Samples were stored at 4 °C. Prior to scanning, samples were allowed to warm to room temperature. NMR spectra were acquired using the first increment of a 1D NOESY pulse sequence with t_{mix} of 100 ms, 4 second acquisition time, 1 s relaxation delay, and a spectral width of 12 ppm.

Following acquisition, spectra were imported into Chenomx NMR Suite 7.5 (Chenomx Inc., Edmonton, Canada) for spectral processing and compound identification and quantification. Phase and baseline corrections were carried out manually referring to the internal standard, DSS, and by

assessing clusters at random across the span of the spectra to make the appropriate adjustments. Shim and chemical shape correction was performed automatically by the software based on a manual comparison of ideal and observed DSS peaks. The known concentration of DSS was used to determine the concentration of all other compounds in the spectrum. Sample pH measurements that were obtained during the sample preparation stage were used for spectrum profiling.

7.2.5.6 Metabolite Profiles

Once spectral processing was completed, compounds were identified and quantified by targeted profiling using Chenomx NMR Suite, 7.0-7.7 (Chenomx Inc., AB). The identification process consisted of manually matching the projection of a predicted compound, selected from the software's built-in 600 MHz compound library, to the processed 1D-¹H NMR spectral signatures. Criteria for placing a compound on the spectrum included both the fit of the projected signature to spectral clusters and background information on metabolite properties, biofunction, pathways and uses. Such metabolite information was obtained from the Chenomx software database, The Human Metabolome Database (HMDB, www.hmdb.ca) and the Kyoto Encyclopedia of Genes and Genomes (KEGG, www.kegg.jp). Quantification of metabolite concentrations were determined by the area of the projected signal once it was fit to the peak centers during identification^{39,41}. As a post hoc evaluation of the metabolite profiles, each compound was assessed based on the number of clusters in the compound's signal, the degree of convolution of the signal with other compounds' signals and the signal to noise ratio. These three characteristics influence the amount of error associated with the identification and quantification process of targeted profiling. Therefore, only compounds that were ranked with high-confidence were included in subsequent analyses of the metabolite profiles.

7.2.5.7 Statistical Analysis

The metabolite profiles were analysed using Principal Component Analysis (PCA) to model the variance between different mixed cultures, and metabolites associated with observed variances. Mean centering ($x - \bar{x}$) and unit variance scaling (x/s , where s is the standard deviation of x concentrations) were applied to the metabolite profiles to easily compare the concentration changes and evaluate all compounds as equally weighted. Metabolite profiles were compared by pair-wise Tukey's Honest Significant Difference (HSD) test where significance was defined by a Bonferroni correction to the alpha value ($p \leq 0.008$). Compounds were identified as a marker when the concentrations in one community's metabolite profile were significantly different from that of all other populations' profiles and where this cross-profile significance only occurred for one community culture. Since the concentrations in the metabolite profiles spanned several orders of magnitude, the compounds within a profile were separated into five bins (20+ mM, 10-20 mM, 1-10 mM, 0.1-1 mM and 0-0.1mM). These bins were defined by a kernel density estimation of compounds' mean concentration, using a bandwidth of 0.031 (Figure A7). The reason for examining the compounds with other compounds of similar magnitude is two-fold. The first is that it allows clear contrast in seeing small increments in concentration changes when visualizing the profiles. The second is that large changes in high-concentration compounds are more likely to be influential than similar percent changes in low-concentration compounds. While this sort of weighting is not desired when analyzing variances and model-building with PCA, it is still important to acknowledge the effect size differences between different compounds when interpreting the changes. All statistical analysis was executed using R, a statistical computing and graphics tool¹⁴⁵. All graphics were generated by the ggplot2 package¹⁴⁶.

7.2.6 Results

7.2.6.1 Multivariate Analysis

In total, 40 metabolite profiles were acquired from the samples collected. PCA was used to visualize samples given the variation in metabolite profiles. This was done via a score plot (Figure 7-2), where each sample is represented by a point on the graph, and its coordinates on the graph demonstrates the relationship of that sample to others given the covariance matrix of the metabolite profiles. Samples from different donors clustered distinctly from each other, and samples of the same donor clustered together. Even the samples of donor A, which were obtained from two different runs, demonstrated similar metabolite profiles and clustered close together. Given that principal component (PC) 1 and 2 cumulatively account for 64.8% of the total variation explained, while PC 3 accounts for only 11.3% of the total variation explained, it can be concluded that donor variability, rather than run-to-run variability, is the main contributor to the amount of variation observed in the metabolite profiles.

Additionally, it became evident that the metabolite profiles derived from human fecal communities were more similar to each other than the profiles of defined communities, despite the latter being designed to resemble the consortia of one of the human donations. Profiles of human feces-derived cultures differed from profiles of the defined culture along PC1 on the x-axis, which accounted for 39.7% of the total variation. The compounds that drove the relationships depicted in the score was extracted from the corresponding loading plot, where the coordinate location of the compounds describes how those compounds contribute to those sample relationships (Figure A8B). The similarities and differences in cultures were also apparent through a global view of the metabolite profiles comparing mean compound concentrations across different donors (Figure 7-3). The compounds most influential in differentiating feces-derived communities and defined mixed cultures

are listed in Table 7-2. Carboxylic acids and derivatives, fatty acids and conjugates, and amino acids and derivatives were the most common types of compounds that differentiated the two groupings.

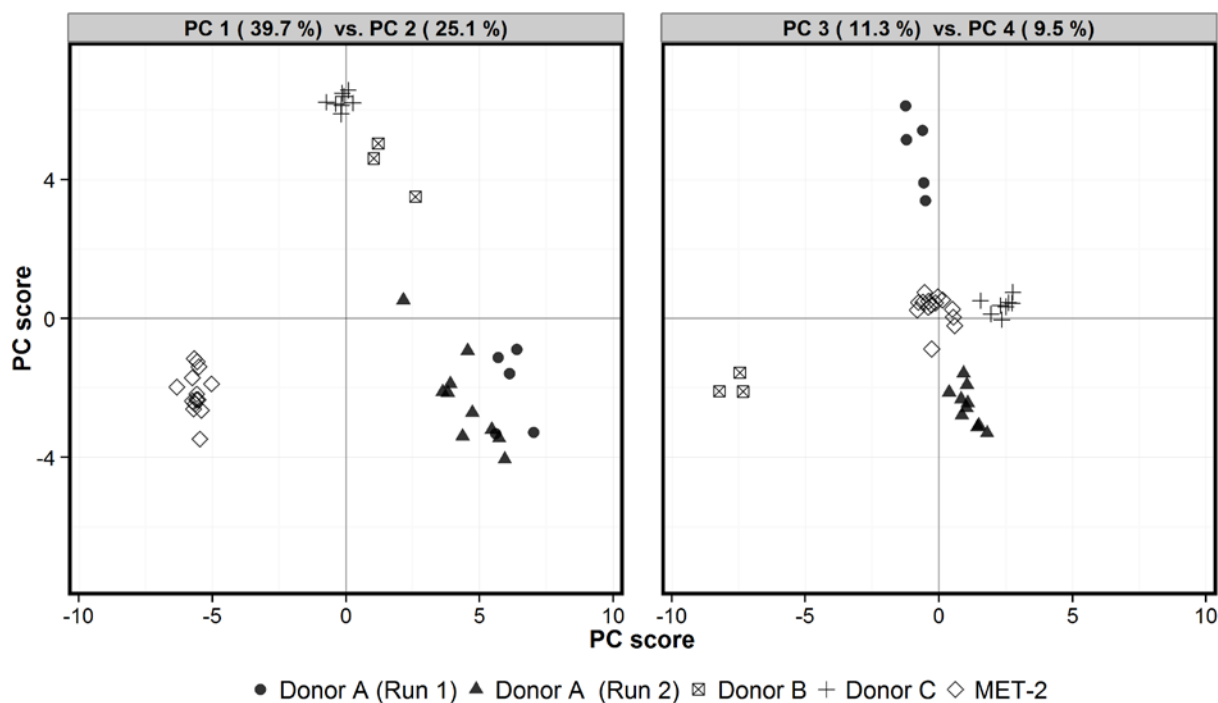


Figure 7-2. Principal component analysis score of metabolite concentration profiles obtained from samples of bioreactor effluent. Metabolite concentrations have been mean centered and unit variance scaled. Seeding culture sources are represented by symbols, which are coloured by runs, where a run includes several bioreactor vessels seeded with the same inoculant and operating in parallel. Principal component (PC) 1 and PC 3 are represented along the x axis, while PC 2 and PC 4 are represented along the y-axis. PC 1 and PC 2 account for 39.7% and 25.1% of the total variation, respectively. PC 3 and PC 4 account for 11.3% and 9.5% of the total variation respectively.

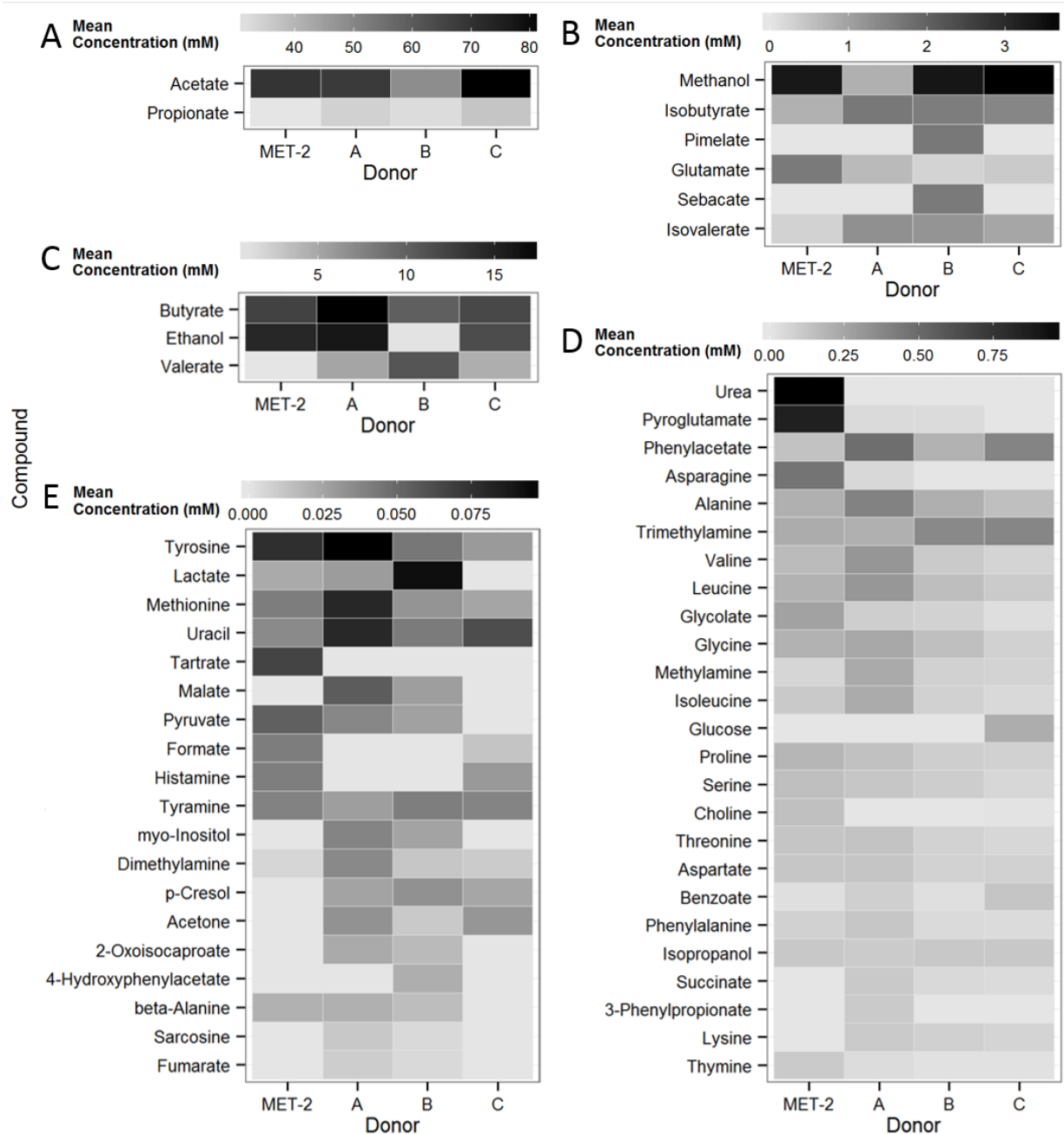


Figure 7-3. Metabolite profile comparison different donors. The mean compound concentrations (mM) of samples from the untreated conditions are presented for each donor. The profiles are split into five bins, as defined by kernel density plot of mean concentrations (Supplementary figure). (A) Bin range: mean ≥ 20 mM. (B) Bin range: $1 \leq \text{mean} < 10$ mM. (C) Bin range: $10 \leq \text{mean} < 20$ mM. (D) Bin range: $0.1 \leq \text{mean} < 1$ mM. (E) Bin range: $0 \leq \text{mean} < 0.1$ mM.

Table 7-2. Compounds that differentiated donor-derived cultures and defined cultures as tested by the Mann-Whitney test at the 95% confidence level.

	Compound⁵	P-value	Compound Class
Bin 1	Propionate	2.76E-06	Carboxylic Acids and Derivatives
Bin 2	Valerate	8.62E-11	Fatty Acids and Conjugates
Bin 3	Isobutyrate	8.62E-11	Carboxylic Acids and Derivatives
	Isovalerate	8.62E-11	Fatty Acids and Conjugates
	Glutamate	2.65E-07	Amino Acids and Derivatives
Bin 4	Phenylacetate	8.62E-11	Phenylacetic Acid Derivatives
	Methylamine	1.72E-10	Alkylamines
	Urea	1.52E-09	Ureas
	Asparagine	3.77E-08	Amino Acids and Derivatives
	Choline	6.70E-08	Alcohols and Polyols
	Succinate	1.44E-07	Carboxylic Acids and Derivatives
	Pyroglutamate	2.38E-07	Pyrrolidines
	Thymine	2.61E-07	Diazines
	Glycolate	2.65E-07	Hydroxy Acids and Derivatives
	Benzoate	3.78E-07	Benzoic Acid and Derivatives
	Lysine	1.10E-05	Amino Acids and Derivatives
	Proline	4.17E-05	Amino Acids and Derivatives
	Serine	0.000155	Amino Acids and Derivatives
	3-Phenylpropionate	0.000655	Cinnamic Acid Derivatives
	Trimethylamine	0.017892	Alkylamines
	Glucose	0.024295	Monosaccharides
	Tartrate	1.52E-09	Carboxylic Acids and Derivatives
	Formate	6.73E-08	Carboxylic Acids and Derivatives
	Histamine	1.10E-07	Azoles
	p-Cresol	1.44E-07	Phenols and Derivatives
	Dimethylamine	3.07E-07	Alkylamines
	Acetone	3.70E-07	Carbonyl Compounds
	2-Oxoisocaproate	9.72E-05	Fatty Acids and Conjugates
	Malate	9.73E-05	Carboxylic Acids and Derivatives
	Fumarate	9.73E-05	Fatty Acids and Conjugates
	myo-Inositol	9.73E-05	Cyclic Alcohols and Derivatives
	Sarcosine	9.73E-05	Amino Acids and Derivatives

⁵Compounds are listed in order of increasing p-value within each compound bin. Bin 1: ≥ 20 mM; Bin 2: ≥ 10 mM, < 20 mM; Bin 3: ≥ 1 mM, < 10 mM; Bin 4: ≥ 0.1 mM, < 1 mM; Bin 5: > 0 mM.

Pyruvate	0.000227	Keto-Acids and Derivatives
Tyramine	0.017207	Phenethylamines
Uracil	0.022443	Diazines
beta-Alanine	0.037663	Amino Acids and Derivatives

7.2.6.2 Metabolite Profiles

Amino acids and fatty acids were the two classes of compounds most commonly profiled in each of the cultures and the abundance of these compounds was a trait common to metabolite profiles of every bacterial community investigated. On average, feces-derived cultures had more compounds in their metabolite profiles compared to that of the defined culture; donor A, B and C had 47, 48 and 40 different compounds respectively and MET-2 had 42 different compounds. When the frequency of each compound class in the metabolite profiles were weighted to the respective concentrations assigned, carboxylic acids and derivatives, fatty acids and conjugates and alcohols and polyols occupied the largest proportions of each culture's profiles (Table 7-3)

Table 7-3. Metabolite profile diversity of each bacterial culture. Proportion of each compound class (%) weighted to the compounds' concentration.

Compound Class	Number of Compounds	Bacterial Culture			
		A	B	C	MET-2
Carboxylic Acids and Derivatives	9	69.59801	74.72572	77.31247	72.48956
Fatty Acids and Conjugates	5	15.95946	19.2173	11.08216	9.650773
Alcohols and Polyols	4	11.324	3.840536	9.901108	13.1407
Amino Acids and Derivatives	18	2.067571	1.280227	0.794359	2.670579
Phenylacetic Acid Derivatives	1	0.323521	0.170689	0.250661	0.09505
Alkylamines	3	0.316235	0.392174	0.295705	0.20722
Hydroxy Acids and Derivatives	2	0.077436	0.142395	0.016511	0.206701
Diazines	2	0.074816	0.046046	0.04838	0.099248
Cinnamic Acid Derivatives	1	0.069176	0	0	0
Benzoic Acid and Derivatives	1	0.055678	0.018187	0.077969	0.016401
Pyrrolidines	1	0.027272	0.031161	0	0.614319
Cyclic Alcohols and Derivatives	1	0.025528	0.022009	0	0
Keto-Acids and Derivatives	1	0.024497	0.022676	0	0.038515
Carbonyl Compounds	1	0.021668	0.00918	0.019653	0
Phenethylamines	1	0.018246	0.035099	0.024367	0.027895
Phenols and Derivatives	2	0.016894	0.046596	0.015592	0
Azoles	1	0	0	0.018856	0.029123
Monosaccharides	1	0	0	0.142208	0
Ureas	1	0	0	0	0.713916

In fact, the proportional distribution of all the compound classes were very similar between all four types of cultures. One of the few discrepancies between the cultures was that donor B-derived culture profiles contained 3.84% alcohols and polyols while all other cultures had ~10% or more of this class of compounds. While characteristics such as compound class composition and number of metabolites profiled could be generalized, each population of enteric bacteria still produced unique metabolite profiles that could be described based on several identifying features.

Each of the communities' metabolite profiles had unique features where certain metabolites could be identified as a marker for that particular bacteria population. 3-Phenylpropionate was unique to the metabolite profiles of donor A-derived cultures as it was not found in any other cultures' metabolite profiles. Donor A-derived samples also had elevated levels of leucine, valine, isoleucine, methylamine, phenylalanine and methionine ($p \leq 0.008$). Mean methanol levels were much lower, while mean succinate, and dimethylamine levels were higher in donor A than all the other cultures, however this was not statistically significant due to relatively large standard deviations.

Donor B culture had three markers that were unique to its metabolite profile. Pimelate and sebacate, both dicarboxylic acids, were found at 1.63 mM and 1.60 mM mean concentrations, respectively, but were absent in the profiles of communities other than Donor B. 4-Hydroxyphenylacetate, an aromatic acid, was found only in donor B culture's profile, though at very low concentrations (mean 0.023 mM). Other identifying features of donor B community's metabolite profile was decreased levels of acetate and ethanol. Donor B cultures had an average concentration of 48.6 mM of acetate and 0.91 mM of ethanol, whereas other communities had significantly greater levels of acetate and ethanol, in the range of 67.62-81.45 mM and 11.5-16.05 mM, respectively. The donor B community also had the highest levels of valerate and lactate of all the communities. The mean concentration of valerate in the metabolite profiles of donor B communities was 10.8 mM,

while the mean of donor A, C and MET-2 communities were 5.02 mM, 4.41 mM and 0.80 mM, respectively. Similarly, the mean lactate concentration in donor B was 0.093 mM, while the mean lactate concentration in donor A and MET-2 communities were 0.028 mM and 0.023 mM respectively.

Donor C was the only community that yielded metabolite profiles that did not contain any lactate. Similarly, donor C culture profiles did not contain any beta-alanine nor pyruvate, while all other cultures' metabolite profiles had mean concentrations within the range of 0.015-0.026 mM and 0.026-0.054 mM respectively. On the other hand, donor C culture's profile was the only one that contained glucose (mean 0.23 mM).

The MET-2 community was the only culture that produced a metabolite profile that contained tartrate and urea. Furthermore, MET-2 culture metabolite profiles contained significantly elevated levels of glutamate, pyroglutamate, asparagine, glycolate, choline, thymine and formate compared to donor A, B and C metabolite profiles ($p \leq 0.008$). Figure 7-4 provides further detail on these identifying markers.

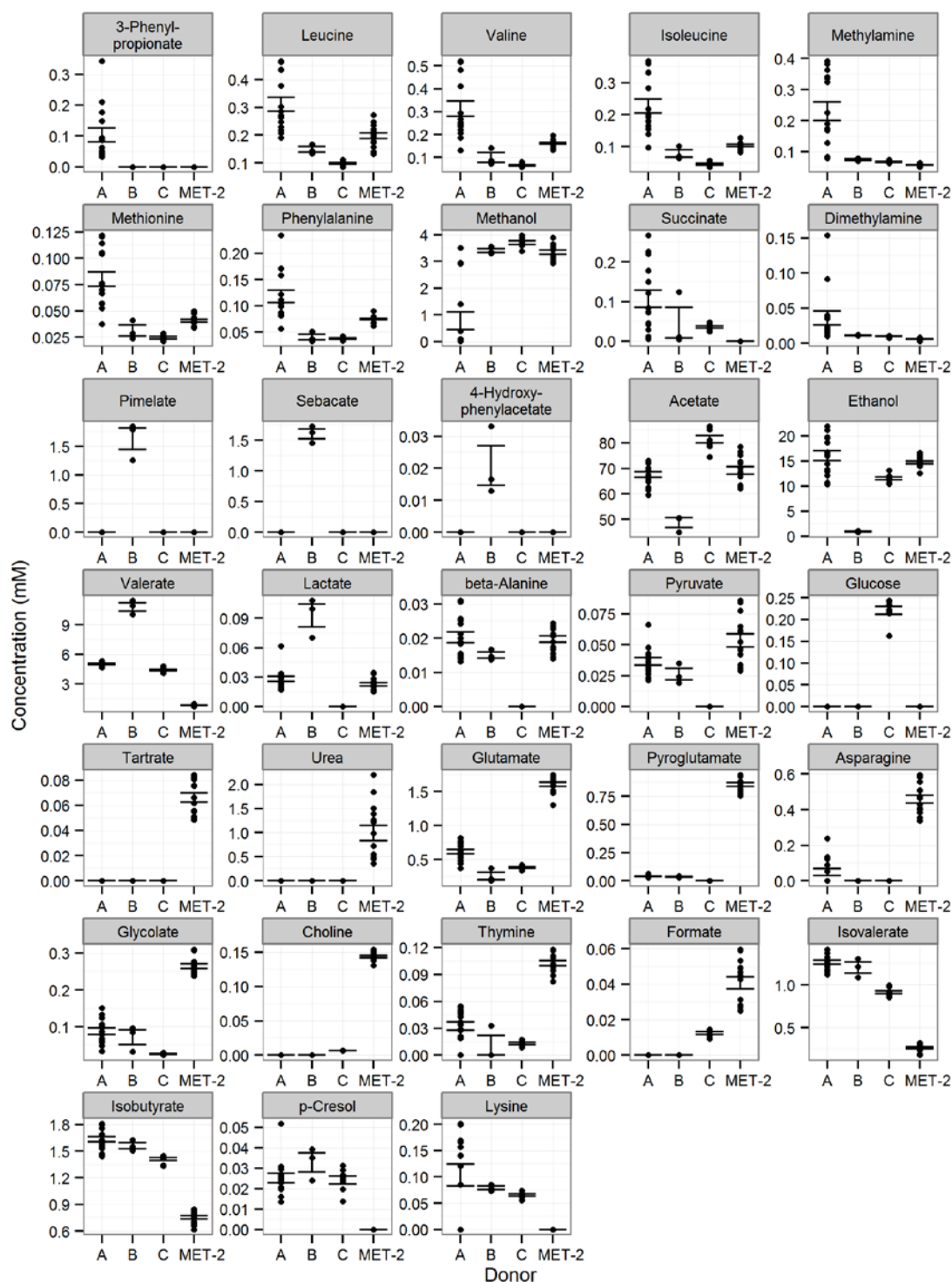


Figure 7-4. Metabolites as identifying markers for metabolite profiles.

7.2.7 Discussion

7.2.7.1 Feces-Derived vs Defined Community Cultures

Metabolite profiles of feces-derived samples were more similar to each other than MET-2 (defined community) samples. Given that the defined community was a subset of isolates derived from a fecal community and which therefore displayed less alpha diversity than its parent community¹⁴⁷, differences in the metabolite profiles between defined and fecal communities were expected. According to multivariate analysis, where compounds are weighted equally, the compounds that were found at greater concentrations in MET-2 compared to fecal cultures were mostly low-concentration compounds (mean < 1 mM) (Figure A8B). The low concentrations of these compounds indicate that these metabolites could be intermediates in major pathways. While these compounds were statistically significant, their impact on the bacterial culture as a whole could be regarded as minimal. Glutamate was the only compound higher in MET-2 culture profiles (mean equal to 1.61 mM) compared to feces-derived culture profiles (means ranged from 0.26 to 0.62 mM). Glutamate is utilized by various species in Groups I through IV of the Clostridia class, forming acetate, butyrate, carbon dioxide, ammonia and hydrogen as products^{148,149}. While acetate and butyrate were observed in abundance, the latter three compounds were not measured in the metabolite profiles due to their incompatibility with ¹H NMR spectroscopy. On the other hand, glutamate synthesis is associated with osmoregulation^{150,151} and ammonia assimilation¹⁵² in anaerobic bacteria. Therefore, the difference in glutamate levels between the feces-derived and defined cultures could be indicative that the communities have differences in their propensity to perform these glutamate-associated biofunctions. Since the biofunctions associated by glutamate synthesis involve fluctuation of intracellular glutamate, those pathways are unlikely to be reflected in the metabolite profiles analyzed here, which represent the exometabolome only. Even though acetate and butyrate were end-products of glutamate

breakdown, they were not observed to be at significantly lower levels in MET-2. This is likely due to the fact that acetate and butyrate are also by-products of other fermentation pathways, such as breakdown of other amino acids¹⁴⁹. Therefore, the concentration of acetate and butyrate observed is the cumulative result of the entire community's metabolic reactions.

Comparison between the two types of cultures also revealed that high-concentration compounds (mean ≥ 1 mM) were found at decreased levels in defined cultures compared to fecal cultures. Since high-concentration compounds were the major components of the metabolite profiles, disparities in those compounds between the two groups are both significant and important. Propionate, valerate, and isovalerate have mean values of 34.7 mM, 6.74 mM, and 1.13 mM respectively in fecal cultures and 30.7 mM, 0.802 mM, and 0.236 mM respectively in the defined culture. Propionate, valerate, isobutyrate, and isovalerate are short chain and branched chain fatty acids (SCFA and BCFA respectively) previously reported to be produced by anaerobic bacteria of the distal GI tract^{74,153,154}. Propionate fermentation occurs via the acrylate pathway in members of the *Clostridium* genus or via the succinate-propionate pathway in other bacteria such as *Bacteroides fragilis*^{155,156}. Since these SCFA and BCFAs were not found in the metabolite profile of the CSTR growth medium, all the propionate, valerate, isobutyrate and isovalerate detected was produced by bacterial metabolic activities.

7.2.7.2 Metabolite Profiles of Bacterial Communities

The metabolite profiles can be broken down into metabolite classes as classified by HMDB (Table 3). The profiles can be viewed as a composition of chemical taxonomy classes proportional to the concentration of each compound. The proportion of compound classes that make up the metabolite profiles was consistent across all the bacterial cultures analyzed, where carboxylic acids and derivatives, fatty acids and conjugates and alcohols and polyols were present in the highest

concentrations, cumulatively accounting for at least 95% of the metabolite profiles. Variability in the metabolite composition occurred in the remaining 2-5% of the profiles. For example, amino acids occupied ~2% of the metabolite profiles in donor A-, and B-derived community and the defined community, but in the donor C-derived community made up only 0.79% of the metabolites. Less than 4% of the metabolites from donor B-derived cultures were alcohols and polyols, but made up ~10% or more of the metabolites in the profiles from other cultures. Despite these small variations, the consistency of the metabolite profiles' structure could be an indication that the intestinal communities share a functional steady state. Therefore, the metabolic interactions tend towards satisfying these functional roles, even when the microbial composition of the community are different.

While the general structure of the metabolite profiles based on the compounds' chemical taxonomic classes could be a characteristic of communities of intestinal bacteria, it does not preclude the fact that there are important metabolic differences between various intestinal microbial ecosystems that are reflected in their metabolite profiles. Such are distinctive and could be used as identifying markers.

7.2.7.3 Donor A

The microbial ecosystem derived from donor A was elevated in compounds involved in the valine-leucine-isoleucine metabolism pathway and the presence of 3-phenylpropionate (absent in all other fecal community samples, including MET-2) suggested an active phenylalanine metabolism pathway^{157,158}. 3-Phenylpropionate is known to be produced by 13 different species in the Clostridiales order, including *C. sporogenes*, *C. caproicum* and *C. botulinum* (*Clostridiaceae* family), and *C. bifermentans* and *C. sordellii*, (*Peptostreptococcaceae* family)^{133,159}, none of which are components of the MET-2 community (Table 1). Alternatively, 3-phenylpropionic acid has been demonstrated to be catabolized by *Escherichia coli* to produce acetaldehyde and pyruvate¹⁵⁷. Given

that MET-2 was designed to mimic the donor A community, and both communities contained *E. coli*, it can be assumed that both donor A and MET-2 communities would have similar, if any, 3-phenylpropionate catabolism capabilities. Therefore, the absence of 3-phenylpropionate in all culture profiles except that of donor A's was likely due to increased phenylalanine dissimilation capabilities, rather than decreased usage of 3-phenylpropionate.

7.2.7.4 Donor B

The metabolite profile of donor B could be distinguished by the presence of pimelate, sebacate and 4-hydroxyphenylacetate. Derivatives of pimelate have been shown to be part of the lysine biosynthesis pathway in organisms such as *E. coli* and *Bacillus subtilis* and the biotin biosynthesis pathway in *B. spaericus*^{160,161}. Bacteria that are able to oxidize long-chain dicarboxylic acid compounds, such as sebacate, are also capable of metabolizing a wide range of dibasic acids, like pimelate, suberate and azelate¹⁶². 4-Hydroxyphenylacetate is an aromatic acid that has been demonstrated to be part of the phenylalanine-tyrosine pathway of *Klebsiella pneumoniae* and the denitrifying pathway of *Pseudomonas* spp.^{163,164}. Acetate and ethanol are produced by fermentation activities of enteric bacteria, and also act as the main substrate for sulfur-reducing enteric bacteria^{154,165}. Given these observations, the donor B-derived community appeared to have greater dicarboxylic metabolism capabilities than the other communities studied, which in turn could be related to either increased levels of fermentation activity or decreased presence of sulfate-reducing bacteria. Valerate and lactate were identified as elevated in donor B-derived communities, and both are associated with fermentation; valerate is a SCFA produced from metabolism of carbohydrates and amino acids¹⁶⁶ and lactate is a hydroxy acid that is produced alongside other organic acids during fermentation^{154,165}.

7.2.7.5 Donor C

Interestingly, neither lactate nor pyruvate were detected in the donor C-derived community metabolite profile, even though these compounds were found in all other communities studied. These are common compounds within microbial communities, however, and we speculate that their concentrations were below the level of detection using our methods, rather than being completely absent. The donor C community metabolite profile was also the only one that still contained glucose. Since complex carbohydrates were the carbon source available in the feed medium, presence of glucose in the effluent samples was indicative of successful breakdown of these polysaccharides but reduced propensity to take up simple sugars. Thus, the absence of lactate and pyruvate as well as the presence of residual glucose all indicated that the donor C microbial community had not metabolized all of the simple sugar. The lack of beta-alanine was also unique to the donor C community. Several organisms, including *C. sporogenes*, follow a reductive pathway in degrading pyrimidines to produce beta-alanine¹⁶⁷, and in some microorganisms beta-alanine can be subsequently broken down to be used as a nitrogen source^{168,169}. As such, the absence of beta-alanine in the donor C profile could be attributed to either a non-reductive pathway of degrading pyrimidines or the inability to utilize beta-alanine.

7.2.7.6 MET-2

MET-2 was the only bacterial population that had tartrate in its metabolite profiles. Tartrate is a dicarboxylic acid known to contribute to fecal pH¹⁷⁰. Certain anaerobic bacteria, including some species found in feces, are able to breakdown tartrate to produce SCFAs, ethanol, lactate, carbon dioxide and hydrogen^{171–174}. The production of tartrate has been studied extensively in *Gluconobacter suboxydans*, a strict aerobe, for industrial purposes^{175,176}. While *G. suboxydans* is not a gut microbe, its tartrate biosynthesis pathway (from glucose, with glycolate as a by-product) could be similar to

those in the MET-2 community. Indeed, glycolate was found in MET-2 profiles at levels an order of magnitude greater than in all other cultures. Based on the evidence that the presence of tartrate was unique to the MET-2 community, and that glycolate was elevated in the MET-2 culture, it appears that the MET-2 community has greater tartrate producing capabilities, rather than a diminished level of tartrate utilization.

Urea was also a unique characteristic of the MET-2 culture. Urea is a waste product from both purine and amino acid degradation, and is usually degraded by bacterial urease to generate free ammonia^{177,178}. Amino acid degradation products were detected in all studied cultures, but urea was only detected in MET-2. We speculate that because MET-2 has low diversity compared to fecal communities, the MET-2 culture lacked the ability to perform urea hydrolysis, rather than had an enhanced ability to break down amino acids.

MET-2 cultures also had higher concentrations of glutamate, pyroglutamate, asparagine, choline, thymine, glycolate and formate. Glutamate and pyroglutamate both participate in glutamate metabolism, such as ammonia assimilation. Asparagine is known to be associated with generation of ammonia via an asparagine-aspartate deamidation pathway^{179,180}. Therefore, an abundance of asparagine could be an indication of decreased activity along the asparagine degradation pathway. Another amino acid that was elevated in MET-2 cultures compared to other communities was choline. In anaerobic bacteria, choline utilization involves its enzymatic cleavage to produce trimethylamine, which can be subsequently used by methanogenic bacteria^{181–183}. Significantly decreased levels of trimethylamine in the MET-2 samples compared to the donor B and C derived communities were most likely because there was a reduced capacity in the MET-2 community to cleave choline. Interestingly, the donor A community, from which MET-2 was derived, also had low levels of trimethylamine, indicating that this characteristic could be preserved in a derivative community.

Thymine is a pyrimidine that is utilized as a nitrogen source via either reductive or oxidative pathways^{184,185}. Given that the MET-2 community exhibited pyrimidine levels over 3 times greater than that of the donor A samples, and 10 times greater than the donor B and C samples, we conclude that MET-2 was not able to catabolize thymine to the same extent as other community cultures. Glycolate and formate were also found to be significantly higher in MET-2 than other cultures. Glycolate, as well as formate, is involved in energy metabolism during fermentation of carbohydrates and degradation of amino acids^{149,186}. Although, as described above, glycolate presence may be attributed to tartrate production. Even though both formate and glycolate were observed at a concentration less than 1 mM in MET-2 samples, this level was still relatively high compared to the profiles of donor A, B or C-derived communities. Subtle differences in metabolically linked compounds, such as glycolate and formate could be representative of differences in the metabolic interactions occurring within each culture.

7.2.8 Conclusion

Communities of colonic bacteria propagated *in vitro* were described by their metabolite profiles, which were built using 1D-¹H NMR spectroscopy. Overall, the metabolite profiles were consistent in the proportions of compound classes, where carboxylic acids (69.6-77.3%), and fatty acids (9.7-19.2%) were the most abundant types of metabolites. The two types of cultures (defined or feces-derived) were seen to have different sugar metabolism capabilities, based on the levels of SCFAs (propionate and valerate) and BCFAs (isovalerate and isobutyrate). The simpler defined community, MET-2, was less proficient at carbohydrate metabolism compared to its feces-derived community counterparts. In general, more compounds were found in the feces-derived cultures, though many that were unique to the feces-derived cultures were found at very low concentrations, less than 0.1 mM.

Ultimately, the features within each metabolic profile not only identified the community culture, but also allowed inferences on their respective metabolic capabilities to be made.

7.2.9 Acknowledgements

This work was supported in part by the generous contribution of Chenomx Inc. in the form of Chenomx NMR Suite 7.0-7.7 software. Acknowledgements are also extended to NSERC Discovery, NSERC Strategic Network (MabNet), and a University of Waterloo's Chronic Disease Prevention Initiative Seed Grant (PROPEL Center for Population Health Impact) Grants to MGA, and an NSERC Canada Graduate Scholarship to SS. Microbiology and CSTR culture work was supported by NSERC Discovery and OMAF/University of Guelph partnership grants to EAV, a CIHR Doctoral Research Award to JAKM, and an Ontario Graduate Scholarship to KS. The authors also thank Michelle Daigneault for her technical expertise in microbial culture and Anita Yen for her artistic contribution to the abstract figure.

7.3 Analysis of Perturbations to Human Fecal Microbiota Propagated *In Vitro*

7.3.1 Experiment Objective

Given that the metabolite profiles were demonstrated to reveal insight into the metabolic functions occurring in *in vitro* microbial communities, further experiments were conducted to assess the robustness of the communities' functional roles. Several different perturbations were administered to feces-derived and defined communities after achieving steady state in the CSTRs. The objective of this study was to characterize the effects of different types of perturbations on bacterial communities based on their metabolite profiles. A wide range of perturbing agents was used to demonstrate the breadth of investigations to which this approach could be used.

7.3.2 Justification of Treatment Condition

7.3.2.1 Clindamycin

Clindamycin is a broad spectrum antibiotic that is active against anaerobic bacteria, such as those found within the gastrointestinal system¹⁸⁷. The antibiotic inhibits protein synthesis within the cell by binding to the 50S ribosomal subunits and preventing aminoacyl sRNA from binding to the messenger ribosome complex¹⁸⁸. This mechanism of action is effective against most anaerobes as well as most aerobic gram-positive bacteria (Table A6), though its ineffectiveness against *Clostridium difficile* cautions clinicians to reserve the use of clindamycin for serious infections^{189,190}. Due to its potency against gram-negative anaerobic bacteria, clindamycin is an appropriate antibiotic to target microflora in the oxygen depleted, distal parts of the GI tract¹⁹⁰.

7.3.2.2 Vancomycin

Vancomycin is a broad spectrum antibiotic usually reserved for treating serious infections caused by β -lactam-resistant pathogens¹⁹¹. Like β -lactam antibiotics, vancomycin targets the mechanisms involved in biosynthesis of bacterial cell walls. While β -lactam antibiotics prevent the crosslinking of peptides to the peptidoglycans of the cell wall, vancomycin operates by specifically targeting the terminal peptides on the peptidoglycans required for transpeptidation reactions in forming the bacterial cell wall^{192,193}.

7.3.2.3 Norepinephrine

Norepinephrine is a catecholamine that functions as a neurotransmitter or a hormone, depending on the area of the body in which the compound is produced. The production of norepinephrine is under control of the sympathetic nervous system, and activates the fight-or-flight response to influences heart rate and other physiological changes that cue the body for rapid metabolic change and physical movement¹⁹⁴. As a hormone, norepinephrine travels through the blood to reach their target tissue to reinforce the sympathetic responses¹⁹⁵. It has been documented that a substantial amount of the body's total norepinephrine production and turnover is performed by mesentery organs, which includes the gastrointestinal tract, spleen and pancreas¹⁹⁶. Since exposure to stress has been shown to induce a sustained increase in catecholamine levels, the effects of psychological stress on the gut microbiota have been probed regarding the integrity of the microflora and resistance to colonization of pathogenic organisms^{197–200}. Specifically, studies in the literature suggest that norepinephrine increases *in vitro* growth of pathogenic gram negative bacteria found in the GI tract: *Yersinia enterocolitica*, *Escherichia coli*, and *Pseudomonas aeruginosa*²⁰¹. It is of interest to investigate if the documented effects of norepinephrine on pure cultures manifest in the same way in community cultures and is able to influence the overall structure of the bioreactor microflora.

7.3.2.4 Defined Bacterial Community

A defined community of bacteria was developed by the Allen-Vercoe group of the University of Guelph as a simplified version of the gut microbiota. Referred to as MET-1, the mixed culture is composed of 33 different strains of bacteria that were isolated from a fecal sample obtained Donor D (Table A1). Using a subset of defined community to perturb the CSTR culture system allows one to monitor the test culture's response to competition.

7.3.3 Materials and Methods

7.3.3.1 Overview of Bioreactor Setup

This study was performed over 6 different bioreactor experimental 'runs', where any given run had one or more bioreactor vessels operating. All vessels within a run propagated cultures from the same seeding source; that is, the inoculant for all vessels in each run was derived from the same inoculum. Conversely, different bioreactor runs used different inoculants. The inocula for Runs 1 and 2 were collected Donor A whereas the inocula for Run 3 was collected from Donor B. Run 5 was inoculated with a defined mixed culture, MET-2, modeled on the fecal community of Donor A (Figure 7-5). Runs 1, 2, and 3 were maintained for 48, 53, and 41 days respectively, while, Run 5 was maintained for 27 days. Based on previous studies using this system the community cultures reached steady-state by ~5 weeks after inoculation⁵.

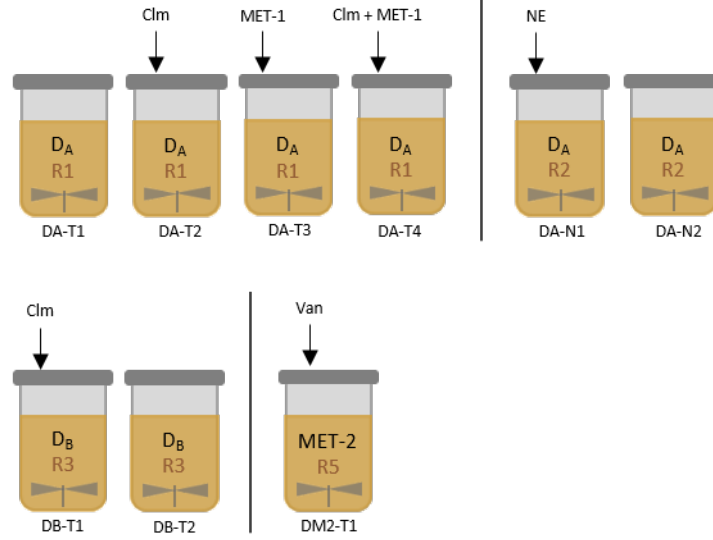


Figure 7-5. Experiment overview. The community culture grown in each bioreactor as labelled as D_A=Donor A, D_B=Donor B, D_C=Donor C, MET 1=Microbial Ecosystem Therapeutic-1, MET-2=Microbial Ecosystem Therapeutic-2, Clm=clindamycin, NE=norepinephrine, Van=vancomycin.

7.3.3.2 Bioreactor Perturbations

One of several perturbing agents were used to treat the mixed cultures by delivering the agent directly to the bioreactor vessel. In Run 2, the Donor-A derived culture (DA-N1) was administered norepinephrine (L-(-)-norepinephrine (+)-bitartrate salt monohydrate, Sigma Aldrich) to a final concentration of 100 mM every 12 hours starting from day 48 post-inoculation until day 52.5 post-inoculation. Similarly, an equal volume of sterile water was added to the sister vessel, DA-N2, as a control. In Run 1, clindamycin (clindamycin 2-phosphate, Sigma Aldrich) was added the Donor A-derived culture (DA-T2) every 12 hours from day 36 to day 40 post-inoculation. The final clindamycin concentration reached at each dosage was 12.6 µg/mL on day 36, 16.6 µg/mL on day 37, 22.4 µg/mL on day 38, 124.7 µg/mL on day 39, and 203.8 µg/mL on day 40. At the same time, the culture DA-T1 was administered dosage-equivalent volumes of sterile water to be assessed as the control vessel. The culture DA-T4 was administered clindamycin from day 36 to day 40 post-inoculation at the same dosages as DA-T2. Following clindamycin treatment, MET-1 was delivered to DA-T4 on day 42 post-inoculation. The MET-1 treatment was prepared as described in Chapter 3, and administered directly into the bioreactor vessel, being careful not to displace the resident bioreactor culture during the treatment delivery process. Concurrently, the DA-T3 community was administered clindamycin-dosage-equivalent volumes of sterile water from day 36-40 and on day 42, was the same MET-1 treatment as DA-T4 was delivered. In Run 3, the culture DB-T1 was administered clindamycin every 12 hours from day 36-40 at the same dosages as those given to DA-T2. As the control vessel, DB-T2 was administered dosage-equivalent volumes every 12 hours from day 36-40. In run 5, the MET-2 culture was treated with vancomycin to a final concentration of 0.5 mg/mL on day 26 post-inoculation.

7.3.3.3 Sample collection and processing

The bioreactor Runs 1, 2, and 3 were studied based on samples from the bioreactor effluent on a daily basis, except for samples from Run 1, where samples were collected from the effluent pooled over a period of 10 days. Runs 4, and 5 were sampled directly from the bioreactor, where sample volumes did not exceed 2.5% of the total working volume so not to impact the chemostatic conditions of the bioreactor. Samples were centrifuged (5000 rpm, 10 min) to remove cells, then subjected to ultracentrifugation (37 000 rpm, 4°C) to pellet media particulates. Finally, samples were filtered using 0.22 µm syringe filters (Acrodisc® Supor membrane, Pall Corporation, UK) to remove bacterial cells. Samples were stored in 4°C until time of NMR analysis.

7.3.3.4 Statistical analysis

The metabolite profiles were analyzed using principal component analysis (PCA) to model the variance between different mixed cultures, and metabolites associated with observed variances. Mean centering ($x - \bar{x}$) and unit variance scaling (x/s , where s is the standard deviation of x concentrations) were applied to the metabolite profiles to easily compare the concentration changes and evaluate all compounds as equally weighted. Changes in specific compounds due to perturbations were assessed using the Mann-Whitney U test, at a 95% confidence level. Since the concentrations in the metabolite profiles spanned several orders of magnitude, the compounds within a profile were separated into five bins (20+ mM, 10-20 mM, 1-10 mM, 0.1-1 mM and 0-0.1mM). These bins were defined by a kernel density estimation of mean concentration of compounds in untreated samples, using a bandwidth of 0.03089 (Figure A7). The reason for examining the compounds with other compounds of similar magnitude is two-fold. The first is that it allows clear contrast in seeing small increments in concentration changes when visualizing the profiles. The second is that large changes in high-concentration compounds are more likely to be influential than similar percent changes in low-

concentration compounds. While this sort of weighting is not desired when analyzing variances and model-building with PCA, it is still important to acknowledge the effects size differences between different compounds when interpreting the changes. All statistical analysis was executed using R, a statistical computing and graphics tool.

7.3.4 Results

7.3.4.1 Antibiotic treatment

Administration of clindamycin and vancomycin to different community cultures produced several changes in the metabolite profiles. Two different fecal communities were treated with clindamycin, one derived from donor A, and one from donor B. The effect of clindamycin treatment was a substantial change in many of the metabolites (Figure 7-6). In the high concentration compounds (mean ≥ 1 mM), which consisted of mostly short chain fatty acids (SCFA), long chain fatty acids (LCFA), branched chain fatty acids (BCFA) and alcohols, there were few drastic changes observed with clindamycin treatment. Most of the fatty acids decreased by a moderate amount, in the 20-50% range, except for acetate, which increased by 17% in the clindamycin-treated community from donor A. Ethanol and glutamate were the only compounds that increased by more than 100% after clindamycin treatment, and both increases were observed in the donor B community. In the Donor A community, the level of ethanol did not change significantly and glutamate increased by 63%. Of the all compounds profiled at lower concentrations (< 1 mM), amino acids and some organic acids were affected the most, experiencing increases of greater than 100%. Succinate, a SCFA, increased by more than 100% after clindamycin treatment. Acetone, a carbonyl, and 2-oxoisocaproate, a fatty acid conjugate, were also low-concentration compounds that experienced large increases. Few low-concentration compounds decreased in concentration with clindamycin treatment, but the ones that did only did so moderately. These included trimethylamine, methylamine, benzoate and

dimethylamine. Some compounds, such as threonine, aspartate, serine, lysine, sarcosine and fumarate, saw significant increases in treatment of Donor B community, but no significant changes in treatment of Donor A community.

The effects of vancomycin were examined by administering the antibiotic to a defined community, MET-2. The overall effects of this treatment were low to moderate decreases in high-concentration compounds: acetate (-11%), propionate (-40%), butyrate (-37%), and isovalerate (-43%). Additionally, there were large increases in several low-concentration compounds: alanine, valine, leucine, glycine, isoleucine, proline, phenylalanine, tyrosine, lactate, methionine, and formate (in order of decreasing mean concentration in untreated samples). Some low-concentration compounds were observed to decrease in concentration after treatment, including trimethylamine (-47%), methylamine (-42%) benzoate (-63%), and tyramine (-33%), and lactate was completely eliminated with vancomycin treatment.

7.3.4.2 Defined community

The effects of challenging the resident community with a defined community were also investigated by administering MET-1, a community culture derived from Donor A. The metabolite profile of the Donor A culture remained largely unchanged after the bacterial community was challenged with MET-1 based on PCA of the metabolite profiles of the culture before and after introduction of the defined culture (Figure A9). However, there were still concentration differences apparent in the donor A profile after MET-1 treatment. High-concentration compounds (mean ≥ 1 mM) either did not undergo significant changes with MET-1 treatment, or decreased by low-to-moderate amounts. The high-concentration compounds that decreased were propionate (-11%), ethanol (-32%), isobutyrate (-12%), glutamate (-18%) and isovalerate (-7%). Methanol was the only high-concentration compound that exhibited a large decrease in concentration (92%) after MET-1 was introduced. The changes

observed in low-concentration compounds were quite varied; observing both moderate increases and decreases in many compounds. Succinate, which was found at less than 1 mM prior to MET-1 treatment, was that only SCFA that increased with MET-1 treatment, and it did so by more than 100%. Similarly, dimethylamine was the only other low-concentration compound that increased by more than 100% after MET-1 treatment. On the other hand, asparagine, lysine and thymine were all completely eliminated after treatment.

The effects of treatment with a defined culture on a bacterial population in dysbiosis were evaluated by treating the Donor A culture with clindamycin followed by adding MET-1 to the culture. The metabolite profile of the culture after clindamycin and MET-1 treatment was comparable to that of cultures treated only with clindamycin according to the PCA score plot, where compounds were assigned equal weighting (Figure A9). In comparing the treatments' profiles without scaling, it was evident that the concentration changes caused by the combination treatment resembled that of the clindamycin treatment, but with some characteristics of the MET-1 treatment profile (Figure 7-6). Most of the changes observed in high-concentration compounds followed the same trend in the clindamycin-MET-1 combination treatment, as clindamycin-only treatment of Donor A. For some compounds, like acetate and glutamate, the percent change was nearly identical for the combination treatment and clindamycin treatment. For other compounds, like butyrate and valerate, the percent differences was more exaggerated in the combination treatment than the clindamycin treatment, with 42% and 73% decrease, respectively, in the combination treatment, compared to 30% and 24% decrease, respectively, in the clindamycin treatment. There were also compounds that followed the trend of MET-1 treatment, rather than the clindamycin treatment, as exemplified by propionate and isovalerate.

Low-concentration compounds (<1 mM) also contained characteristics from both clindamycin treatment and MET-1 treatment profiles. In fact, about one third of the low-concentration compounds followed the same trends as the changes seen with clindamycin treatment and one-third followed that of MET-1 treatment. The remaining third of the compounds exhibited changes that were unique to the combination treatment, resembling neither the clindamycin treatment nor the MET-1 treatment. Specifically, valine, leucine, and isoleucine did not change significantly when clindamycin was followed with introduction of MET-1 even though those compounds increased with both clindamycin-only and MET-1-only treatments. Glycolate (100+%), serine (40%), threonine (72%) and aspartate (86%) all increased after the combination treatment, but demonstrated no significant changes after clindamycin treatment, and moderate decreases with MET-1 treatment. 3-Phenylpropionate, lactate and malate were also compounds that demonstrated concentration changes after administration of the combination treatment that were different from the changes observed when treated with each constituent perturbation agent alone.

7.3.4.3 Norepinephrine treatment

Exposure of norepinephrine, a stress hormone, to the Donor A-derived culture had a minor impact on the culture's metabolite profiles, as reflected in the PCA of treated and non-treated profiles (Figure A9). In an overview of changes occurring across the entire profile (Figure 7-6), there were very few significant changes, giving further support that little changes occurred due to norepinephrine treatment. The few changes that did occur were observed in low-concentration compounds (mean \leq 1 mM). These changes occurred in succinate (61%), malate (77%), 3-phenylpropionate (100+%), fumarate (100+%) and asparagine (-100%).

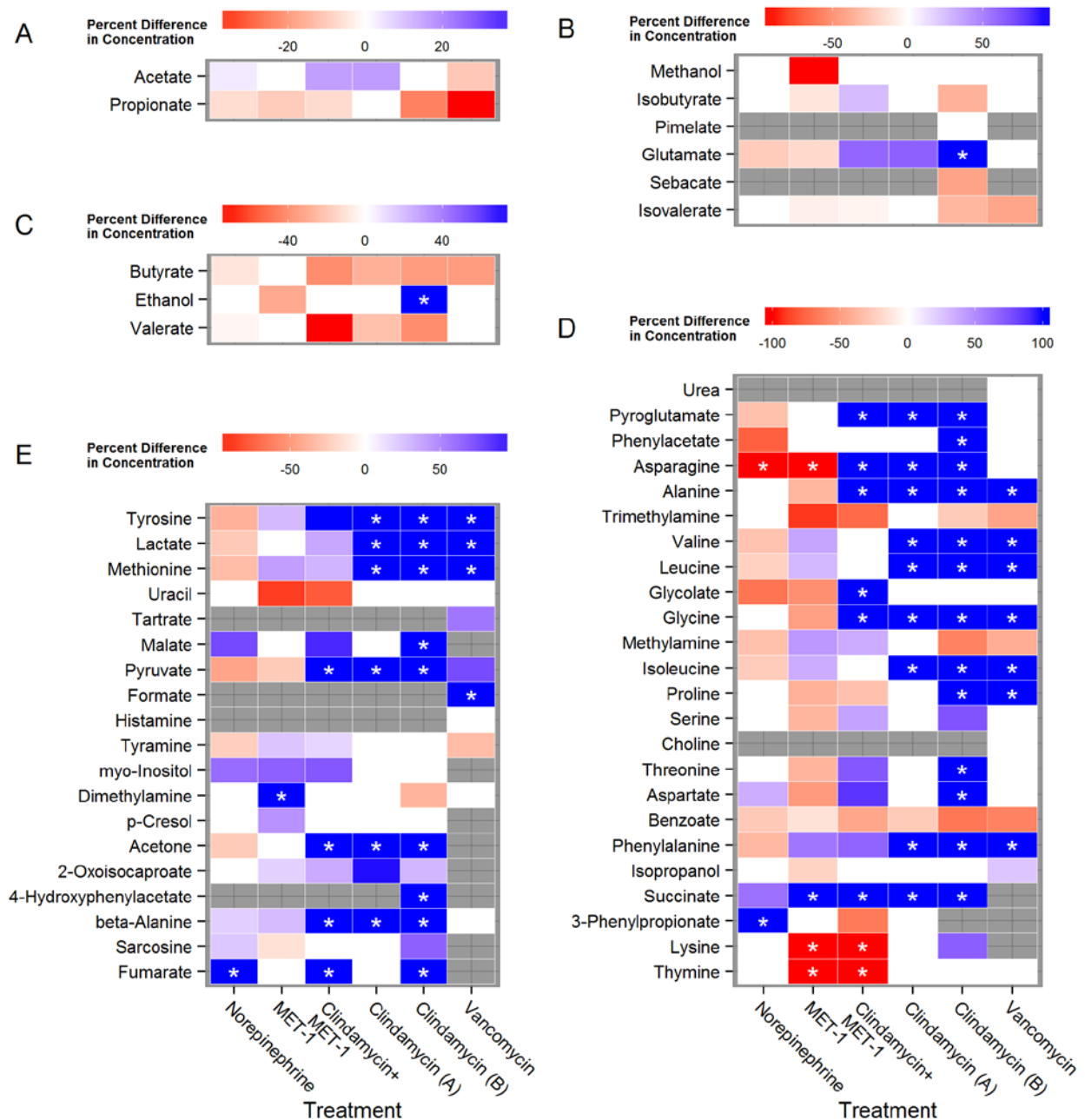


Figure 7-6. Treatment effects on metabolite profiles. The percent difference between mean compound concentration of treated samples and mean compound concentration of their non-perturbed counterparts. Norepinephrine, MET-2, clindamycin-MET-2 combination treatment and clindamycin treatment were administered to bioreactors inoculated with Donor A. Clindamycin treatment was also administered to bioreactors inoculated with Donor B. The two clindamycin conditions are identified by their respective donors as Clindamycin (A) and Clindamycin (B). Vancomycin was administered to a bioreactor growing the defined culture, MET-2. Positive percent differences are represented in blue, negative percent differences are represented in red, and non-significant differences are

represented in white. Compounds not found in either treatment or non-treated conditions for a particular profile are presented grey. Differences exceeding 100% are represented by a star (*) on blue, while differences exceeding -100% are represented by a star on red. Significance of compound changes due to treatment was tested by the Mann-Whitney U Test, where p-value > 0.05 were considered as non-significant. Compounds are organized according to the mean of untreated samples. (A) Bin range: mean \geq 20 mM. (B) Bin range: 1 mM \leq mean < 10 mM. (C) Bin range: 10 mM \leq mean < 20 mM. (D) Bin range: 0.1 mM \leq mean < 1 mM. (E) Bin range: 0 mM \leq mean < 0.1 mM.

7.3.5 Discussion

7.3.5.1 Clindamycin treatment

The impact of clindamycin on the metabolic output of both Donor A and Donor B communities was fairly consistent, with percent changes mirrored in almost every metabolite. The majority of compounds that increased by more than 100% after clindamycin treatment were amino acids. This indicates that pathways involving the metabolism of amino acids were strongly impacted by clindamycin treatment. Indeed, many changes in compound concentrations correspond with known metabolic pathways of anaerobic bacteria. For example, alanine and asparagine all increased more than 100% after clindamycin treatment; glutamate increased by 62.7% in Donor A and more than 100% in Donor B. In *B. fragilis*, the alanine-aspartate-glutamate metabolism pathway is responsible for the breakdown and formation of each of those amino acids²⁰². More generally, compounds involved in the glycine-serine-threonine metabolism increased after clindamycin treatment. Glycine, and methionine both increased by more than 100% after antibiotic treatment. Threonine and serine increased significantly in clindamycin treatment of Donor B culture, but no effect was observed in Donor A culture. Valine-leucine-isoleucine biosynthesis was another pathway where all three end products of the pathway increased by more than 100%. Aromatic amino acids such as phenylalanine, and tyrosine also increased. Given that clindamycin acts on the protein synthesis mechanisms of bacteria to inactivate most anaerobic and Gram-positive aerobic bacteria^{187,188,203}, the increased amounts of free amino acids observed further supports that extensive impairment of amino acid metabolic pathways had occurred following clindamycin administration.

Short chain fatty acids metabolism pathways were also affected by clindamycin treatment. This is apparent in the decreased levels of propionate, butyrate, valerate, isobutyrate and isovalerate. This supports findings from previous reports of decreases in fatty acid output levels of fecal bacteria

in response to clindamycin treatment²⁰⁴. Production of short-chain fatty acids occurs due to carbohydrate fermentation by clostridia. Since, in general, members of the genus *Clostridium*, with the exception of *C. difficile*, are susceptible to clindamycin, decreased levels of short-chain fatty acids could be attributed to elimination of organisms that perform carbohydrate fermentation. Interestingly, succinate increased by more than 100% in both donor communities treated with clindamycin. Since succinate is a SCFA closely associated with energy metabolism pathways as well as carbohydrate fermentation in anaerobic bacteria, the large change in succinate levels could be a reflection of drastic shifts in the communities' energy usage^{205,206}.

7.3.5.2 Vancomycin treatment

Similar to clindamycin treatment, administration of vancomycin caused a reduction in SCFA: propionate, butyrate, and isovalerate. Other SCFAs observed in the metabolite profiles, valerate and isobutyrate, did not change significantly with treatment. This indicated that the antibiotic inactivated the organisms that were metabolizing carbohydrates and producing SCFAs, such as vancomycin-sensitive bacteria from the *Clostridium leptum* and *Clostridium coccoides* clusters²⁰⁷. SCFAs, and lactate are fecal metabolites that have been previously reported to decrease with vancomycin treatment in mice models and human subjects^{208,209}, even though the findings presented here demonstrated a large increase in lactate. The valine-leucine-isoleucine metabolism pathway as well as phenylalanine-tyrosine-tryptophan pathway were both well represented in the list of amino acids that increased after antibiotic treatment. These amino acid changes are in agreement with previous publications on the metabolomic effects of vancomycin on gut microbiota²⁰⁸. Glycine, serine and threonine, which are part of the same metabolism pathway, all increased after vancomycin treatment. Glycine is also involved in energy storage by acting as an oxidizing agent, and it is reduced to acetate and ammonia in one-carbon metabolism pathways^{149,210}.

7.3.5.3 Challenge of Fecal Culture with Defined Community

The impact of using a defined mixed culture, MET-1, to challenge an already-established fecal culture was investigated. In comparing fecal culture from Donor A before and after it was challenged with MET-1, the metabolite profiles were not impacted strongly by the presence of the defined culture. There were low to moderate decreases in SCFAs that had mean concentrations ≥ 1 mM prior to treatment. This was an indication that the Donor A community's ability to ferment carbohydrates was not strongly impacted by the introduction of MET-1. Even though succinate increased by over 100%, such small concentrations were observed that the 100+% increase still only amounted to less than 1 mM. Despite from the small magnitude of change, the doubling of succinate should not be discounted as it could be an indication that MET-1 treatment had implications on Donor A community's energy consumption or conservation pathways²⁰⁵. Methanol did not change significantly in all other treatment conditions tested, whereas there was a 91% decrease of methanol with MET-1 treatment. While this was an observation unique to MET-1 treatment, the absolute change in methanol was only 0.7 mM, indicating that the overall impact on the metabolite profile, and the culture as a whole, was very little. Other metabolite-specific changes can be observed with introduction of MET-1, however, like methanol, the effect size of those changes were minimal.

7.3.5.4 Challenge of Fecal Culture in Dysbiosis with Defined Community

MET-1 was also introduced to Donor A fecal culture after it was treated with clindamycin. This was driven by the hypothesis that the defined community would be able to restore a culture in dysbiosis back to the original culture. The metabolite profiles revealed that while MET-1 was not able re-establish the culture back to its original community, the profile of the clindamycin-MET-1 treated culture contained characteristics of the clindamycin-treated Donor A community, the MET-1 treated Donor A community as well as some characteristics unique to the combination treatment. Similar

trends were observed in short-chain fatty acids such as acetate, butyrate, and valerate, indicating that impairment to carbohydrate fermentation pathways by clindamycin persisted even after introduction of MET-1. The increase in amino acids associated with the alanine-glutamate-aspartate pathway and phenylalanine-tyrosine-tryptophan pathway increased, just like the clindamycin-treated cultures. Therefore, certain amino acid degradation pathways impacted by clindamycin treatment were not recovered by the introduction of MET-1.

On the other hand, common features between the MET-1 treatment profile and the combination treatment profile was evidence that some of the metabolic functions of MET-1 were established. For example, decreases in metabolically linked compounds trimethylamine, and methylamine^{211,212} observed in both MET-1-treated communities could mean that the pathway involved in the of these compounds became more active with the introduction of MET-1. Similarly, thymine and uracil, both catabolized by *E. coli* as carbon sources in nitrogen-restricted environments, decreased after MET-1 treatment and combination treatment²¹³. Some amino acid metabolism in the combination-treated profile was more like the MET-1-treated profile than the clindamycin treated profiles. This included decreases in proline, and lysine, and an increases in methionine. The similar trends that proline and lysine follow in both MET-1 treatment and combination treatment condition suggested that MET-1 established its ability to breakdown nitrogen-rich amino acids. Even though free nitrogen and ammonia levels were not measured, the increase in catabolism of pyrimidines and nitrogen-rich amino acids were evidence that MET-1 conferred its nitrogen-utilisation capabilities to the Donor A community, even when it was perturbed with clindamycin.

As mentioned earlier, the combination-treatment condition displayed some unique characteristics that were not observed in either MET-1 treatment or the combination treatment conditions. In particular, amino acids in the leucine-valine-isoleucine pathway remained the same

after combination treatment, whereas clindamycin-treated only communities exhibited large increases and the MET-1-treated community exhibited moderate increases in those three compounds. The opposing effects of the combination treatment compared to the clindamycin treatment suggested that MET-1 had a remedial effect with respects to the degradation pathway of these amino acids. Furthermore, compounds of the glycine-serine-threonine pathway seemed to be sensitive to difference types of perturbations. Increases in serine and threonine was unique to the combination treatment; methionine increased moderately after both MET-1 treatment and the combination treatment; glycine increased by more than 100% in both combination treatment and clindamycin treatment, but not MET-1 treatment. 3-Phenylprionate decreased by 67% after the combination treatment, an effect not observed with either MET-1 treatment or clindamycin treatment. Since the Donor A community was previously found to harbour phenylalanine degradation capabilities, the decrease in 3-phenylpropionate suggested that the combination treatment impaired the pathway responsible for converting phenylalanine to 3-phenylpropionate, an effect that did not occur with single perturbation treatments. Compounds closely associated with anaerobic energy metabolism, like glycolate, fumarate, lactate, malate and succinate^{186,205}, also responded to the combination treatment differently than the MET-1 or clindamycin single-treatments. Glycolate, malate and fumarate all increased largely after the combination treatment, but either did not change significantly, or decreased moderately in the single-treatment conditions.

Based on direct comparison of the combination-treated metabolite profiles with other treatment, the signatures of clindamycin treatment profiles were clearly evident. At the same time, it was apparent that MET-1 was able to confer some of its metabolic functions. Interestingly, some of the effects of clindamycin treatment was no longer evident after the combination treatment, suggesting that the MET-1 treatment re-established particular biofunctions that were impaired by the

antibiotic. Since the metabolite profile of the Donor A community after the combination treatment appeared to be a mosaic of clindamycin treatment and MET-1 treatment, it suggested that the community was undergoing a dynamic, transitional stage, where the MET-1 community was in the process of establishing its functional role in the “deteriorated” Donor A community. It was speculated that if the MET-1 community was given more time to assert itself following clindamycin treatment, the metabolite profile of the combination-treated community would become distinct from either clindamycin or MET-1 treatment.

The changes highlighted in the heat maps demonstrated how small changes in the metabolite profiles can yield insight into the changes happening to the community. At the same time, the impact of these changes on the profile as a whole is minimal. As a result, the metabolite profiles of Donor A community after combination treatment and after clindamycin treatment can still be considered to be very similar.

7.3.5.5 Norepinephrine treatment

Previous studies have shown that physiological stress can alter gut microbiota¹⁹⁸. Norepinephrine, a catecholamine released in response to acute stress, is known to promote the growth of, and expression of virulence determinants by, some Gram-negative enteric bacteria, such as *E. coli* and *Y. enterocolitica*²⁰¹, through bacterial signalling pathways. However, based on multivariate analysis the metabolic profiles of a norepinephrine-treated community did not differ notably from the control (untreated) community. A difference in succinate was one of the few changes seen in response to norepinephrine treatment, although this change was not as drastic in comparison to succinate fluctuations seen in response to antibiotic perturbation. One mechanism by which norepinephrine enhances growth of certain enteric organisms involves liberation of glycoprotein-bound iron ions by the catecholamine, rendering iron available to the bacteria to be taken up, promoting their growth. It

is possible that the addition of norepinephrine to the bioreactor system did not have an effect on microbial community growth because iron was supplied in the form of hemin, which is readily taken up by anaerobic bacteria^{214,215}, and not trapped in a protein-iron complex.

Chapter 8 Conclusions

8.1 Metabolomic analysis of Nutrient Solutions

The targeted profiling technique used to build metabolite profiles of nutrient solutions was able to identify and quantify compounds based on their peak signals in ^1H NMR spectra. Despite the manual spectral processing and peak assignment, the technique has been validated to be one that is reproducible and reliable^{39,41}. Conducting metabolomic analysis on cell culture media provides insight into the nutrient content of the growth solution. Therefore, as a demonstration of the metabolomic analysis in the context of nutrient solutions, the metabolite make up of a defined media for mammalian CHO cells subjected to various treatments of UV irradiation was investigated. The findings from this study revealed that not only was the constructed metabolite profiles a comprehensive annotation of the nutrient solution, the profiles were also sensitive enough to quantify the minimal effects of UV treatment. Furthermore, this experiment demonstrated that the metabolite profiling technique would be an appropriate tool in observing changes in low-molecular weight compounds found in cell culture media. Therefore, the application of this technique can be expanded to investigating metabolites of spent media, from the growth of bacterial communities, such as those of the human gut microbiota.

8.2 Standardization of Reporting on Fecal Water Metabolites

Metabolomic studies of the human gut microbiota thus far have focused heavily on fecal samples. However, the minimal degree of standardization of obtaining NMR-based metabolomic data of feces⁶² and the complete lack of standardization in reporting the fecal metabolomic data have prevented cross-experimental comparisons. Therefore, with the purpose of determining an appropriate means of reporting fecal metabolome, various concentrations of human fecal water were analyzed. NMR-based analysis of these fecal water samples revealed that there was a direct linear relationship

between fecal concentration and metabolite concentration. The implication of this is that the concentration of biomass used in preparing such samples must be well documented and transparent to all experiment stakeholders, or be reported as normalized values. Two methods of normalizing metabolite profiles were investigated: one where the concentrations were scaled to the mass of feces from which the metabolites were extracted, and the other, where the concentrations were scaled to an internal compound that was observed in abundance. While both normalization methods were equally effective, they both had advantages and disadvantages in their application. Furthermore, given the linear relationship between metabolite concentration and fecal concentration, the most suitable fecal-water concentration would be defined by the intended analysis method. For the purposes of NMR-based metabolite concentration, fecal concentration should be greater than 25% fecal concentration in order to ensure low concentration compounds are captured within the limit of quantification. Alternatively, it can be expressed that fecal water samples must contain at least ~7.3 mM of acetate to construct accurate and encompassing metabolite profiles. Beyond analytical considerations, practical limitations of sample preparation should also be taken into account. Since all samples must be passed through a 0.22 μ m filter, this may serve as an upper limit to the concentration of fecal water that can be processed.

Based on these findings, it is recommended that acetate be used as a proxy to qualifying the contents of a fecal water sample. Acetate is an attractive identifying marker because it has been consistently found to be the most abundant compound in fecal water samples. According to the samples evaluated here, fecal water samples should be evaluated to contain at least 7.3 mM or 0.43 g/L of acetate in order to qualify for metabolomic analysis by 1D- ^1H NMR. Describing a sample in terms of acetate content could be a more standardized method of describing fecal water samples. While 7.3 mM of acetate was the minimum requirement determined from this experiment, further

study on other stool samples would be required to determine a generalized acetate threshold for fecal water samples.

8.3 Metabolomic Evaluation of *In Vitro* Gut Bacterial Communities

A wealth of insight has been gained on the human gut microbiota through the use of *in vitro* growth of the gut microbes as a community culture. These *in vitro* systems have been validated for their ability to preserve the *in vivo* gut bacterial community on the microbiological level^{11,68,75}, but little work has been put forth to determine the metabolomic implications of these *in vitro* gut microbial cultures. In an effort to provide such metabolomic descriptions of *in vitro* gut bacterial communities, two experiments were conducted. The first centered on developing an appropriate procedure for preparing bioreactor samples for NMR-based metabolomic analysis. The second investigated the metabolomic impact of the microbial make-up of the fecal inoculant on its resulting *in vitro* bacterial community.

8.3.1 Sample Preparation Procedure for Metabolomic Analysis

Comparison of various procedures for the preparation of bioreactor samples for NMR-based metabolomic analysis determined that factors such as performing a “soft spin” to remove cell matter and syringe filtration had no significant effect on the metabolite profiles. On the other hand, the centrifugal force impacted the metabolite profile. It was found that increasing centrifugal force applied the effluent samples generally decreased compound concentrations in a linear fashion. These linear decreases, while significant, only occurred within 2% of the mean compound concentration and can be considered to be of little importance. The only compound to increase in concentration with centrifugal force applied was p-cresol, which increased by approximately 50% per 10 000 rpm applied. Furthermore, these effects were found to occur with high reproducibility. In conclusion,

centrifugal force had significant, but minimal impact on the metabolite profile of the sample, but additional attention should be allocated to p-cresol. The large and predictable means by which this compound is impacted by ultracentrifugation means that a correction factor of 50% per 10 000 rpm should be applied to the compound's concentration. Compensation for other compounds is not recommended as their trends were not observed at as high precision, and changes made to their concentrations due to sample processing steps are so small as to be considered negligible.

8.3.2 Effects of Varying Inoculant Biomass Proportions

There were four different community cultures evaluated in this experiment: two different defined cultures, MET-2A and MET-3A and their respective reciprocal cultures, MET-2B and MET-3B. The reciprocal cultures consisted of the same strains as the parent culture (MET-2A or MET-3A), but with different proportions of each strain in the inoculum. The objective of this experiment was to observe the metabolite profiles of the culture over time in a CSTR system and determine if similar metabolic states would be achieved regardless of inoculant biomass. After determining sources of variation to be mostly attributed to different inocula, rather than experimental parameters such as run variability or profiling variability, the effect of varying inoculum biomass was investigated. First of all, the total concentration of metabolites in a defined culture was similar to the total metabolite concentration in the profiles of the respective reciprocal population. Evaluation of the cultures over time revealed that the majority of compounds appeared to approach a common level by day 20 when comparing the MET-2A and MET-3A with their respective reciprocal cultures. In order to draw more accurate conclusions of the trends in metabolite profile changes over time, however, metabolite profiles should be assembled for more time points throughout the duration of the culture, with multiple samples per time point for an estimate of error and variation and also for a longer duration in order to evaluate achievement of metabolic steady state.

8.4 Metabolomic Characterization of *In Vitro* Gut Microbial Communities

The advantages associated with studying the gut microbiota in an *in vitro* system combined with the metabolomic data available in these systems allows for a wide range of studies that contribute to metabolomic characterization the human gut microbiota. Therefore, two experiments were conducted with the aim of providing further insight into the metabolome of these bacterial communities. In “Metabolomic Analysis of Human Fecal Microbiota Propagated *In Vitro* in a CSTR System,” it was found that the metabolite profiles of the gut bacterial communities were generally composed of the same proportion of classifications of metabolites, regardless of the individual from whom the inoculating community was obtained. Despite these general similarities, the metabolite profiles were still able to distinguish one bacterial community from the next, even identifying differences in metabolic capabilities. These metabolomic characterizations were also applied to a defined community cultured in the same *in vitro* bioreactor system, revealing some functional differences between feces-derived and defined culture communities.

Given that the metabolite profiles were able to identify individual bacterial communities and comment on their metabolic functions, it was of interest to observe the *in vitro* communities’ responses to various perturbations. For this purpose, several perturbation agents were administered to feces-derived or defined communities and metabolite profiles were used to identify the changes undergone by the cultures. These treatments included two antibiotic treatments (clindamycin and vancomycin); challenging the resident community with a defined community, MET-2; an antibiotic-MET-2 combination treatment; and a norepinephrine treatment. It was found that amino acid and carbohydrate metabolism were the most sensitive to perturbation, as various changes in levels of associated amino acids and fatty acids were observed in all treatments except for norepinephrine treatment, which did not drastically impact the community’s metabolite profile. Succinate was a

notable compound as it increased considerably after every treatment, even norepinephrine treatment. Due to succinate's close association with central metabolic processes and is subject to large changes with perturbation, it could be the first indicator of a perturbation to the culture. From this study, it has been demonstrated that metabolite profiles of mixed cultures allows the differentiation of communities of different inoculant sources. These metabolite profiles also allow characterization of various perturbations, such as antibiotics, probiotic treatment and stress exposure, on those communities.

The complex symbiotic relationship between microbial ecosystems residing in the human GI tract and human health emphasizes the importance of developing a comprehensive understanding of the interactions through alternative and complementary approaches. Studying the metabolome of the gut microbes is one such approach that has the capability of deciphering those interactions. In this body of work, an *in vitro* system of studying the human gut microbiota was subjected to metabolomic evaluation. Standardizing reports on fecal metabolomes and determining source of variation in metabolomic analyses are fundamental characterizations that appropriately contextualizes the findings presented here so they can be effectively communicated to other researchers. The metabolomic approach was further applied in analyzing the exometabolome found in the nutrient solutions in which these community cultures were grown. The metabolite profiles defined microbial metabolite products that commented on the summative outcome of the complex microbial interactions occurring within the gut microbiota. This is relevant to clinical and industrial research because this metabolomic characterization has the potential to bridge microbial activities with human health and physiology.

Chapter 9 Recommendations

The insights ascertained in this thesis contribute to the metabolomic understanding of the human gut microbiota by setting foundational work from which this approach can be further developed. For example, further metabolomic evaluation of the single-stage CSTR system can be performed. Comparison of the *in vitro* communities derived from two inocula that differed only in biomass proportions of the same bacterial species demonstrated that communities grown in the bioreactor systems were relatively robust, with the final ecological consortium driven by metabolic function, rather than the proportional composition of the inoculant. This finding suggested that the metabolic steady state was being approached as of 20 days post-inoculation. Some fluctuation in the metabolite levels indicated that the system was still dynamic, but convergence of profiles between a culture (MET-2A or MET-3A) and its reciprocal variant (MET-2B or MET-3B) was a hint that metabolic interactions were becoming more stable. Therefore, it would be worthwhile to conduct an experiment in which a metabolomic time-trend can establish metabolic steady state is achieved and how this compares to when ecological steady state is reached.

There is an evident gap in human gut microbial research, where advanced *in vitro* systems have been developed to study the ecology of the human gut microbiota, and at the same time, metabolomic understanding of the human gut microbiota have focused on fecal metabolites. A direct comparison of the fecal metabolome with the metabolome of the *in vitro* community cultured from the same fecal sample would be a good depiction of how the fecal communities change or are preserved in the *in vitro* system. This analysis would be a starting point for bridging that gap between *in vitro* gut microbiota research and our metabolomic understanding of that ecosystem.

The metabolite profile analysis of feces-derived communities and defined communities was able to confer a wealth of information regarding the gut microbiota. Not only did it highlight gross

differences between each community, such as decreased levels carbohydrate fermentation in the defined community compared to its fecal counterpart, but specialized metabolic functions were identified to be unique to certain communities. Since the metabolomic data was able to divulge details on the metabolic interactions in the bacterial consortia, this approach can be applied to studying the human gut microbiota in clinical or pharmaceutical contexts. A wide range of perturbations were administered to the gut as a demonstration of such applications. From this preliminary experiment, it was demonstrated that amino acid metabolism and carbohydrate fermentation were the most dramatically influenced by antibiotic treatments. It is worth noting that the antibiotics used were administered to different kinds of bacterial communities (feces-derived or defined). Given that clindamycin and vancomycin targeted cells through different mechanisms of action – clindamycin inhibits the ribosomal mechanisms to halt protein interaction, while vancomycin interrupts biosynthesis of bacterial cell walls – it would be interesting to see how different classes of antibiotics impact one bacterial consortium.

Furthermore, challenging a feces-derived community with a defined community and the combination treatment of clindamycin followed by challenging with a defined community also yielded interesting results. While the clindamycin-only treatment and the defined community-only treatment both produced unique metabolomic responses in the bacterial community, the combination treatment elicited a metabolite profile that resembled a hybrid of the two stand-alone treatments. This experiment would benefit from extending the time during which the microbial responses are monitored. In observing the *in vitro* cultures until steady-state is reached again, the full, long-term impact of the treatments, whether it be stand-alone or combination treatments, on the metabolic and ecologic makeup of the community would be captured. Additionally, the effect of recovery time on the final outcome in the combination treatment should be investigated. Such a study would reveal

whether the stability of the community impacts the ability of the challenging consortia to establish itself within the residing community.

The norepinephrine treatment was an example of using the *in vitro* communities to investigate host physiology on the gut microbiota. Even though host stress response have been show to impact bacterial physiology and the health of gut microbiota as a whole, those interactions were not reproduced with the when the *in vitro* feces-derived community was treated with norepinephrine. The reason for this discrepancy was hypothesized to be attributed the abundance of iron supplied in the media feed. Since iron was readily available throughout the culturing process, norepinephrine did not have any added benefit to the microbes that would typically capitalize on the norepinephrine-facilitated release on iron ions was lost.

The work presented in this thesis provides a kernel of insight into the contributions metabolomic studies could have on understanding the human gut microbiota. The recommendations stemming from this work focus on validating the bioreactor systems involved in propagating gut microbes *in vitro* as well as in deciphering metabolic interactions of those communities.

References

- (1) Mapelli, V., Olsson, L., and Nielsen, J. (2008) Metabolic footprinting in microbiology: methods and applications in functional genomics and biotechnology. *Trends Biotechnol.* 26, 490–7.
- (2) Saric, J., Wang, Y., Li, J., and Coen, M. (2007) Species variation in the fecal metabolome gives insight into differential gastrointestinal function. *J. Proteome Res.* 7, 352–360.
- (3) Le Gall, G., Noor, S. O., Ridgway, K., Scovell, L., Jamieson, C., Johnson, I. T., Colquhoun, I. J., Kemsley, E. K., and Narbad, A. (2011) Metabolomics of fecal extracts detects altered metabolic activity of gut microbiota in ulcerative colitis and irritable bowel syndrome. *J. Proteome Res.* 10, 4208–18.
- (4) Jacobsen, N. (2007) NMR spectroscopy explained: simplified theory, applications and examples for organic chemistry and structural biology, pp 353–450.
- (5) McDonald, J. a K., Schroeter, K., Fuentes, S., Heikamp-Dejong, I., Khursigara, C. M., de Vos, W. M., and Allen-Vercoe, E. (2013) Evaluation of microbial community reproducibility, stability and composition in a human distal gut chemostat model. *J. Microbiol. Methods* 95, 167–74.
- (6) Kim, B.-S., Kim, J. N., and Cerniglia, C. E. (2011) In vitro culture conditions for maintaining a complex population of human gastrointestinal tract microbiota. *J. Biomed. Biotechnol.* 2011, 838040.
- (7) Fiehn, O., Kopka, J., Dörmann, P., Altmann, T., Trethewey, R. N., and Willmitzer, L. (2000) Metabolite profiling for plant functional genomics 1157–1161.
- (8) Nicholson, J. K., Lindon, J. C., and Holmes, E. (1999) “Metabonomics”: understanding the metabolic responses of living systems to pathophysiological stimuli via multivariate statistical analysis of biological NMR spectroscopic data. *Xenobiotica.* 29, 1181–9.
- (9) Phelps, T. J., Palumbo, A. V, and Beliaev, A. S. (2002) Metabolomics and microarrays for improved understanding of phenotypic characteristics controlled by both genomics and environmental constraints. *Curr. Opin. Biotechnol.* 13, 20–24.
- (10) Fiehn, O. (2002) Metabolomics--the link between genotypes and phenotypes. *Plant Mol. Biol.* 48, 155–71.
- (11) Madsen, R., Lundstedt, T., and Trygg, J. (2010) Chemometrics in metabolomics--a review in human disease diagnosis. *Anal. Chim. Acta* 659, 23–33.

- (12) Gibney, M. J., Walsh, M., Brennan, L., Roche, H. M., German, B., and van Ommen, B. (2005) Metabolomics in human nutrition: opportunities and challenges. *Am. J. Clin. Nutr.* 82, 497–503.
- (13) Mashego, M. R., Rumbold, K., De Mey, M., Vandamme, E., Soetaert, W., and Heijnen, J. J. (2007) Microbial metabolomics: past, present and future methodologies. *Biotechnol. Lett.* 29, 1–16.
- (14) Vigneau-Callahan, K. E., Shestopalov, A. I., Milbury, P. E., Matson, W. R., and Kristal, B. S. (2001) Characterization of diet-dependent metabolic serotypes: analytical and biological variability issues in rats. *J. ...* 131, 924S–932S.
- (15) Gross, J. H. (2014) Mass Spectrometry 2nd ed., pp 355–406. Springer, Heidelberg.
- (16) Takahashi, H., Kai, K., Shinbo, Y., Tanaka, K., Ohta, D., Oshima, T., Altaf-Ul-Amin, M., Kurokawa, K., Ogasawara, N., and Kanaya, S. (2008) Metabolomics approach for determining growth-specific metabolites based on Fourier transform ion cyclotron resonance mass spectrometry. *Anal. Bioanal. Chem.* 391, 2769–82.
- (17) Bristow, A. W. T., Webb, K. S., Lubben, A. T., and Halket, J. (2004) Reproducible product-ion tandem mass spectra on various liquid chromatography/mass spectrometry instruments for the development of spectral libraries. *Rapid Commun. Mass Spectrom.* 18, 1447–54.
- (18) De Vos, R. C. H., Moco, S., Lommen, A., Keurentjes, J. J. B., Bino, R. J., and Hall, R. D. (2007) Untargeted large-scale plant metabolomics using liquid chromatography coupled to mass spectrometry. *Nat. Protoc.* 2, 778–91.
- (19) Lisec, J., Schauer, N., Kopka, J., Willmitzer, L., and Fernie, A. R. (2006) Gas chromatography mass spectrometry-based metabolite profiling in plants. *Nat. Protoc.* 1, 387–96.
- (20) Gao, X., Pujos-Guillot, E., Martin, J.-F., Galan, P., Juste, C., Jia, W., and Sebedio, J.-L. (2009) Metabolite analysis of human fecal water by gas chromatography/mass spectrometry with ethyl chloroformate derivatization. *Anal. Biochem.* 393, 163–75.
- (21) Gao, X., Pujos-Guillot, E., and Sébédio, J.-L. (2010) Development of a quantitative metabolomic approach to study clinical human fecal water metabolome based on trimethylsilylation derivatization and GC/MS analysis. *Anal. Chem.* 82, 6447–56.
- (22) Dumas, M.-E., Maibaum, E. C., Teague, C., Ueshima, H., Zhou, B., Lindon, J. C., Nicholson, J. K., Stamler, J., Elliott, P., Chan, Q., and Holmes, E. (2006) Assessment of

analytical reproducibility of ^1H NMR spectroscopy based metabonomics for large-scale epidemiological research: the INTERMAP Study. *Anal. Chem.* 78, 2199–208.

(23) Beckonert, O., Keun, H. C., Ebbels, T. M. D., Bundy, J., Holmes, E., Lindon, J. C., and Nicholson, J. K. (2007) Metabolic profiling, metabolomic and metabonomic procedures for NMR spectroscopy of urine, plasma, serum and tissue extracts. *Nat. Protoc.* 2, 2692–703.

(24) Wevers, R. A., Engelke, U., and Heerschap, A. (1994) High-Resolution ^1H NMR Spectroscopy of Blood Plasma for Metabolic Studies 40, 1245–1250.

(25) Monakhova, Y. B., Schäfer, H., Humpfer, E., Spraul, M., Kuballa, T., and Lachenmeier, D. W. (2011) Application of automated eightfold suppression of water and ethanol signals in ^1H NMR to provide sensitivity for analyzing alcoholic beverages. *Magn. Reson. Chem.* 49, 734–9.

(26) Hwang, T., and Shaka, A. (1995) Water suppression that works. Excitation sculpting using arbitrary wave-forms and pulsed-field gradients. *J. Magn. Reson. Ser. A* 112, 275–279.

(27) Viant, M. R. (2003) Improved methods for the acquisition and interpretation of NMR metabolomic data. *Biochem. Biophys. Res. Commun.* 310, 943–948.

(28) Tang, H., Wang, Y., Nicholson, J. K., and Lindon, J. C. (2004) Use of relaxation-edited one-dimensional and two dimensional nuclear magnetic resonance spectroscopy to improve detection of small metabolites in blood plasma. *Anal. Biochem.* 325, 260–272.

(29) Chang, D., Banack, C. D., and Shah, S. L. (2007) Robust baseline correction algorithm for signal dense NMR spectra. *J. Magn. Reson.* 187, 288–92.

(30) Jacob, D., Deborde, C., and Moing, A. (2013) An efficient spectra processing method for metabolite identification from (1)H-NMR metabolomics data. *Anal. Bioanal. Chem.*

(31) Weljie, A. M., Newton, J., Mercier, P., Carlson, E., and Slupsky, C. M. (2006) Targeted profiling: quantitative analysis of ^1H NMR metabolomics data. *Anal. Chem.* 78, 4430–42.

(32) Reynolds, W., and Enriquez, R. (2002) Choosing the best pulse sequences, acquisition parameters, postacquisition processing strategies, and probes for natural product structure elucidation by NMR. *J. Nat. Prod.* 65, 221–244.

(33) Beckwith-Hall, B. M., Nicholson, J. K., Nicholls, A. W., Foxall, P. J. D., Lindon, J. C., Connor, S. C., Abdi, M., Connelly, J., and Holmes, E. (1998) Nuclear Magnetic Resonance Spectroscopic and Principal Components Analysis Investigations into Biochemical Effects of Three Model Hepatotoxins. *Chem. Res. Toxicol.* 11, 260–272.

- (34) Holmes, E., and Antti, H. (2002) Chemometric contributions to the evolution of metabonomics: mathematical solutions to characterising and interpreting complex biological NMR spectra. *Analyst* 127, 1549–1557.
- (35) Penha, R., and Hines, J. (2001) Using principal component analysis modeling to monitor temperature sensors in a nuclear research reactor. *Proc. Maint. Reliab.* ... 6–9.
- (36) Jolliffe, I. T. (2002) Principal Component Analysis, Second Edition 2nd ed., pp 1–61. Springer-Verlag New York, Inc.
- (37) Duarte, I. F., Lamego, I., Rocha, C., and Gil, A. M. (2009) NMR metabonomics for mammalian cell metabolism studies. *Bioanalysis* 1, 1597–614.
- (38) Lindon, J., and Nicholson, J. (2000) Metabonomics: metabolic processes studied by NMR spectroscopy of biofluids. *Concepts Magn. Reson.* 12, 289–320.
- (39) Sokolenko, S., Blondeel, E., Azlah, N., George, B., Schulze, S., Chang, D., and Aucoin, M. G. (2014) Profiling convoluted single-dimension proton NMR spectra: A Plackett–Burman approach for assessing quantification error of metabolites in complex mixtures with application to cell culture. *Anal. Chem.* 86, 3330–3337.
- (40) Tredwell, G. D., Behrends, V., Geier, F. M., Liebeke, M., and Bundy, J. G. (2011) Between-person comparison of metabolite fitting for NMR-based quantitative metabolomics. *Anal. Chem.* 83, 8683–7.
- (41) Sokolenko, S., McKay, R., Blondeel, E. J. M., Lewis, M. J., Chang, D., George, B., and Aucoin, M. G. (2013) Understanding the variability of compound quantification from targeted profiling metabolomics of 1D-1H-NMR spectra in synthetic mixtures and urine with additional insights on choice of pulse sequences and robotic sampling. *Metabolomics* 9, 887–903.
- (42) Kaderbhai, N. N., Broadhurst, D. I., Ellis, D. I., Goodacre, R., and Kell, D. B. (2003) Functional genomics via metabolic footprinting: monitoring metabolite secretion by *Escherichia coli* tryptophan metabolism mutants using FT-IR and direct injection electrospray mass spectrometry. *Comp. Funct. Genomics* 4, 376–91.
- (43) Piper, P. (1995) The heat shock and ethanol stress responses of yeast exhibit extensive similarity and functional overlap. *FEMS Microbiol. Lett.* 134, 121–127.
- (44) Devantier, R., Scheithauer, B., Villas-Bôas, S. G., Pedersen, S., and Olsson, L. (2005) Metabolite profiling for analysis of yeast stress response during very high gravity ethanol fermentations. *Biotechnol. Bioeng.* 90, 703–14.

- (45) Aranibar, N., Borys, M., Mackin, N. a, Ly, V., Abu-Absi, N., Abu-Absi, S., Niemitz, M., Schilling, B., Li, Z. J., Brock, B., Russell, R. J., Tymiak, A., and Reily, M. D. (2011) NMR-based metabolomics of mammalian cell and tissue cultures. *J. Biomol. NMR* 49, 195–206.
- (46) Bradley, S. a, Ouyang, A., Purdie, J., Smitka, T. a, Wang, T., and Kaerner, A. (2010) Fermentanomics: monitoring mammalian cell cultures with NMR spectroscopy. *J. Am. Chem. Soc.* 132, 9531–3.
- (47) Read, E. K., Bradley, S. a, Smitka, T. a, Agarabi, C. D., Lute, S. C., and Brorson, K. a. (2013) Fermentanomics informed amino acid supplementation of an antibody producing mammalian cell culture. *Biotechnol. Prog.* 29, 745–53.
- (48) Khetan, A., Huang, Y., Dolnikova, J., Pederson, N. E., Wen, D., Yusuf-Makagiansar, H., Chen, P., and Ryll, T. (2010) Control of misincorporation of serine for asparagine during antibody production using CHO cells. *Biotechnol. Bioeng.* 107, 116–23.
- (49) Qin, J., Li, R., Raes, J., Arumugam, M., Burgdorf, K. S., Manichanh, C., Nielsen, T., Pons, N., Levenez, F., Yamada, T., Mende, D. R., Li, J., Xu, J., Li, S., Li, D., Cao, J., Wang, B., Liang, H., Zheng, H., Xie, Y., Tap, J., Lepage, P., Bertalan, M., Batto, J.-M., Hansen, T., Le Paslier, D., Linneberg, A., Nielsen, H. B., Pelletier, E., Renault, P., Sicheritz-Ponten, T., Turner, K., Zhu, H., Yu, C., Li, S., Jian, M., Zhou, Y., Li, Y., Zhang, X., Li, S., Qin, N., Yang, H., Wang, J., Brunak, S., Doré, J., Guarner, F., Kristiansen, K., Pedersen, O., Parkhill, J., Weissenbach, J., Bork, P., Ehrlich, S. D., and Wang, J. (2010) A human gut microbial gene catalogue established by metagenomic sequencing. *Nature* 464, 59–65.
- (50) Human, T., and Project, M. (2012) Structure, function and diversity of the healthy human microbiome. *Nature* 486, 207–14.
- (51) Arumugam, M., Raes, J., Pelletier, E., Le Paslier, D., Yamada, T., Mende, D. R., Fernandes, G. R., Tap, J., Bruls, T., Batto, J.-M., Bertalan, M., Borruel, N., Casellas, F., Fernandez, L., Gautier, L., Hansen, T., Hattori, M., Hayashi, T., Kleerebezem, M., Kurokawa, K., Leclerc, M., Levenez, F., Manichanh, C., Nielsen, H. B., Nielsen, T., Pons, N., Poulain, J., Qin, J., Sicheritz-Ponten, T., Tims, S., Torrents, D., Ugarte, E., Zoetendal, E. G., Wang, J., Guarner, F., Pedersen, O., de Vos, W. M., Brunak, S., Doré, J., Antolín, M., Artiguenave, F., Blottiere, H. M., Almeida, M., Brechot, C., Cara, C., Chervaux, C., Cultrone, A., Delorme, C., Denariáz, G., Dervyn, R., Foerstner, K. U., Friss, C., van de Guchte, M., Guedon, E., Haimet, F., Huber, W., van Hylckama-Vlieg, J., Jamet, A., Juste, C., Kaci, G., Knol, J., Lakhdari, O., Layec, S., Le Roux, K., Maguin, E., Mérieux, A., Melo Minardi, R., M’rini, C., Muller, J., Oozeer, R., Parkhill, J., Renault, P., Rescigno, M., Sanchez, N., Sunagawa, S., Torrejon, A., Turner, K., Vandemeulebrouck, G., Varela, E., Winogradsky, Y., Zeller, G., Weissenbach, J., Ehrlich, S. D., and Bork, P. (2011) Enterotypes of the human gut microbiome. *Nature* 473, 174–80.

(52) Yatsunenkov, T., Rey, F. E., Manary, M. J., Trehan, I., Dominguez-Bello, M. G., Contreras, M., Magris, M., Hidalgo, G., Baldassano, R. N., Anokhin, A. P., Heath, A. C., Warner, B., Reeder, J., Kuczynski, J., Caporaso, J. G., Lozupone, C. a, Lauber, C., Clemente, J. C., Knights, D., Knight, R., and Gordon, J. I. (2012) Human gut microbiome viewed across age and geography. *Nature* 486, 222–7.

(53) Li, M., Wang, B., Zhang, M., Rantalainen, M., Wang, S., Zhou, H., Zhang, Y., Shen, J., Pang, X., Zhang, M., Wei, H., Chen, Y., Lu, H., Zuo, J., Su, M., Qiu, Y., Jia, W., Xiao, C., Smith, L. M., Yang, S., Holmes, E., Tang, H., Zhao, G., Nicholson, J. K., Li, L., and Zhao, L. (2008) Symbiotic gut microbes modulate human metabolic phenotypes. *Proc. Natl. Acad. Sci. U. S. A.* 105, 2117–22.

(54) Martin, F., and Sprenger, N. (2010) Dietary modulation of gut functional ecology studied by fecal metabonomics. *J. Proteome Res.* 9, 5284–5295.

(55) Jacobs, D. M., Deltimple, N., van Velzen, E., van Dorsten, F. a, Bingham, M., Vaughan, E. E., and van Duynhoven, J. (2008) (1)H NMR metabolite profiling of feces as a tool to assess the impact of nutrition on the human microbiome. *NMR Biomed.* 21, 615–26.

(56) Shinohara, K., Ohashi, Y., Kawasumi, K., Terada, A., and Fujisawa, T. (2010) Effect of apple intake on fecal microbiota and metabolites in humans. *Anaerobe* 16, 510–5.

(57) Cross, A. J., Greetham, H. L., Pollock, J. R. A., Rowland, I. R., and Bingham, S. A. (2006) Variability in fecal water genotoxicity, determined using the comet assay, is independent of endogenous N-nitroso compound formation attributed to red meat consumption. *Environ. Mol. Mutagen.* 47, 179–184.

(58) Jansson, J., Willing, B., Lucio, M., Fekete, A., Dicksved, J., Halfvarson, J., Tysk, C., and Schmitt-Kopplin, P. (2009) Metabolomics reveals metabolic biomarkers of Crohn's disease. *PLoS One* 4, e6386.

(59) Monleón, D., Morales, J. M., Barrasa, A., López, J. A., Vázquez, C., and Celda, B. (2009) Metabolite profiling of fecal water extracts from human colorectal cancer. *NMR Biomed.* 22, 342–8.

(60) Ponnusamy, K., Choi, J. N., Kim, J., Lee, S.-Y., and Lee, C. H. (2011) Microbial community and metabolomic comparison of irritable bowel syndrome faeces. *J. Med. Microbiol.* 60, 817–27.

(61) Ikeda, N., Saito, Y., Shimizu, J., Ochi, a, Mizutani, J., and Watabe, J. (1994) Variations in concentrations of bacterial metabolites, enzyme activities, moisture, pH and bacterial composition between and within individuals in faeces of seven healthy adults. *J. Appl. Bacteriol.* 77, 185–94.

- (62) Wu, J., An, Y., Yao, J., Wang, Y., and Tang, H. (2010) An optimised sample preparation method for NMR-based faecal metabonomic analysis. *Analyst* 135, 1023–30.
- (63) Zheng, X., Xie, G., Zhao, A., Zhao, L., Yao, C., Chiu, N. H. L., Zhou, Z., Bao, Y., Jia, W., Nicholson, J. K., and Jia, W. (2011) The footprints of gut microbial-mammalian co-metabolism. *J. Proteome Res.* 10, 5512–22.
- (64) Antunes, L. C. M., Han, J., Ferreira, R. B. R., Lolić, P., Borchers, C. H., and Finlay, B. B. (2011) Effect of antibiotic treatment on the intestinal metabolome. *Antimicrob. Agents Chemother.* 55, 1494–503.
- (65) Martin, F.-P. J., Wang, Y., Sprenger, N., Yap, I. K. S., Lundstedt, T., Lek, P., Rezzi, S., Ramadan, Z., van Bladeren, P., Fay, L. B., Kochhar, S., Lindon, J. C., Holmes, E., and Nicholson, J. K. (2008) Probiotic modulation of symbiotic gut microbial-host metabolic interactions in a humanized microbiome mouse model. *Mol. Syst. Biol.* 4, 157.
- (66) Wang, X., and Gibson, G. (1993) Effects of the in vitro fermentation of oligofructose and inulin by bacteria growing in the human large intestine. *J. Appl. Bacteriol.* 75, 373–380.
- (67) Aura, A.-M., Oikarinen, S., Mutanen, M., Heinonen, S.-M., Adlercreutz, H. C. T., Virtanen, H., and Poutanen, K. S. (2006) Suitability of a batch in vitro fermentation model using human faecal microbiota for prediction of conversion of flaxseed lignans to enterolactone with reference to an in vivo rat model. *Eur. J. Nutr.* 45, 45–51.
- (68) Campbell, W. L., Franklin, W., and Cerniglia, C. E. (1992) Validation studies on an in vitro semicontinuous culture system designed to simulate a bacterial ecosystem of the human intestine. *J. Microbiol. Methods* 16, 239–252.
- (69) Allison, C., Mcfarlan, C., and Macfarlane, G. T. (1989) Studies on mixed populations of human intestinal bacteria grown in single-stage and multistage continuous culture systems. *Appl. Environ. Microbiol.* 55, 672–678.
- (70) Gibson, G. R., and Roberfroid, M. B. (1995) Dietary modulation of the human colonic microbiota: introducing the concept of prebiotics. *J. Nutr.* 125, 1401–1412.
- (71) Fuller, R. (1989) Probiotics in man and animals. *J. Appl. Bacteriol.* 66, 365–378.
- (72) Duncan, S., Scott, K., Ramsay, A. G., Harmsen, H. J., Welling, G. W., Stewart, C. S., and Flint, H. J. (2003) Effects of alternative dietary substrates on competition between human colonic bacteria in an anaerobic fermentor system. *Appl. Environ. Microbiol.* 69, 1139–1142.

- (73) Steer, T. E., Johnson, I. T., Gee, J. M., and Gibson, G. R. (2007) Metabolism of the soyabean isoflavone glycoside genistin in vitro by human gut bacteria and the effect of prebiotics. *Br. J. Nutr.* 90, 635.
- (74) Bearne, C. a, Mallett, a K., Rowland, I. R., and Brennan-Craddock, W. E. (1990) Continuous culture of human faecal bacteria as an in vitro model for the colonic microflora. *Toxicol. In Vitro* 4, 522–5.
- (75) Macfarlane, G., Macfarlane, S., and Gibson, G. (1998) Validation of a Three-Stage Compound Continuous Culture System for Investigating the Effect of Retention Time on the Ecology and Metabolism of Bacteria in the Human Colon. *Microb. Ecol.* 35, 180–7.
- (76) Molly, K., Woestyne, M., and Verstraete, W. (1993) Development of a 5-step multi-chamber reactor as a simulation of the human intestinal microbial ecosystem. *Appl. Microbiol. Biotechnol.* 1, 254–258.
- (77) Nakamura, N., Lin, H. C., McSweeney, C. S., Mackie, R. I., and Gaskins, H. R. (2010) Mechanisms of microbial hydrogen disposal in the human colon and implications for health and disease. *Annu. Rev. Food Sci. Technol.* 1, 363–95.
- (78) Feria-Gervasio, D., Denis, S., Alric, M., and Brugère, J.-F. (2011) In vitro maintenance of a human proximal colon microbiota using the continuous fermentation system P-ECSIM. *Appl. Microbiol. Biotechnol.* 91, 1425–33.
- (79) Feria-Gervasio, D., Tottey, W., Gaci, N., Alric, M., Cardot, J.-M., Peyret, P., Martin, J.-F., Pujos, E., Sébédio, J.-L., and Brugère, J.-F. (2014) Three-stage continuous culture system with a self-generated anaerobia to study the regionalized metabolism of the human gut microbiota. *J. Microbiol. Methods* 96, 111–8.
- (80) Minekus, M., Smeets-Peeters, M., Bernalier, a., Marol-Bonnin, S., Havenaar, R., Marteau, P., Alric, M., Fonty, G., and Huis in't Veld, J. H. J. (1999) A computer-controlled system to simulate conditions of the large intestine with peristaltic mixing, water absorption and absorption of fermentation products. *Appl. Microbiol. Biotechnol.* 53, 108–114.
- (81) Minekus, M., Marteau, P., Havenaarl, R., and Huis, J. H. J. (1995) A Multicompartmental Dynamic Computer-controlled Model Simulating the Stomach and Small Intestine. *Altern. to Lab. Anim.* 23, 197–209.
- (82) Rajilić-Stojanović, M., Maathuis, A., Heilig, H. G. H. J., Venema, K., de Vos, W. M., and Smidt, H. (2010) Evaluating the microbial diversity of an in vitro model of the human large intestine by phylogenetic microarray analysis. *Microbiology* 156, 3270–81.

- (83) Van Nuenen, M. H. M. C., Diederick Meyer, P., and Venema, K. (2003) The effect of various inulins and *Clostridium difficile* on the metabolic activity of the human colonic microbiota in vitro. *Microb. Ecol. Health Dis.* 15, 137–144.
- (84) Venema, K., van Nuenen, M. H. M. C., van den Heuvel, E. G., Pool, W., and van der Vossen, J. M. B. M. (2003) The effect of lactulose on the composition of the intestinal microbiota and short-chain fatty acid production in human volunteers and a computer-controlled model of the proximal large intestine. *Microb. Ecol. Health Dis.* 15, 94–105.
- (85) Maathuis, A. J. H., Heuvel, E. G. Van Den, Schoterman, M. H. C., and Venema, K. (2012) Galacto-oligosaccharides have prebiotic activity in a dynamic in vitro colon model using a ¹³C-labeling technique. *J. Nutr.* 142, 1205–1212.
- (86) Egert, M., de Graaf, A. a, Maathuis, A., de Waard, P., Plugge, C. M., Smidt, H., Deutz, N. E. P., Dijkema, C., de Vos, W. M., and Venema, K. (2007) Identification of glucose-fermenting bacteria present in an in vitro model of the human intestine by RNA-stable isotope probing. *FEMS Microbiol. Ecol.* 60, 126–35.
- (87) Child, M. W., Kennedy, A., Walker, A. W., Bahrami, B., Macfarlane, S., and Macfarlane, G. T. (2006) Studies on the effect of system retention time on bacterial populations colonizing a three-stage continuous culture model of the human large gut using FISH techniques. *FEMS Microbiol. Ecol.* 55, 299–310.
- (88) To, J., Blunt, B., and Chang, D. (2014) Investigating the criteria for accurate quantitative results with Chenomx NMR Suite, pp 1–4.
- (89) Ihaka, R., and Gentleman, R. (1996) R: a language for data analysis and graphics. *J. Comput. Graph. Stat.* 5, 299–314.
- (90) McDonald, J. A. K. (2013) Development of an in vitro fermentation model to culture the human distal gut microbiota. University of Guelph.
- (91) Schroeter, K. F. B. (2014) Adaptation of a Single-Stage Chemostat System to Model Functional Aspects of the Human Distal Gut Microbiota. University of Guelph.
- (92) De Boever, P., Deplancke, B., and Verstraete, W. (2000) Fermentation by gut microbiota cultured in a simulator of the human intestinal microbial ecosystem is improved by supplementing a soygerm powder. *J. Nutr.* 130, 2599–606.
- (93) Baylis, S. A., Finsterbusch, T., Bannert, N., Blümel, J., and Mankertz, A. (2011) Analysis of porcine circovirus type 1 detected in Rotarix vaccine. *Vaccine* 29, 690–7.

- (94) Aranha, H. (2011) Virus safety of biopharmaceuticals: Absence of evidence is not evidence of absence. *Contract Pharma*.
- (95) Kerr, A., and Nims, R. (2010) Adventitious viruses detected in biopharmaceutical bulk harvest samples over a 10 Year Period. *PDA J. Pharm. Sci. Technol.* 64, 481–485.
- (96) Denyer, S., and Baird, R. (1990) Guide to microbiological control in pharmaceuticals (Rubinstein, M. H., Ed.) 1st ed., pp 158–240. Ellis Horwood Limited, Chichester.
- (97) Yaman, A. (2012) Methods of sterilization for controlled release injectable implantable preparations, in *Long Acting Injections and Implants* (Wright, J. C., and Burgess, D. J., Eds.), pp 459–473. Springer US, Boston, MA.
- (98) Mosley, G. (2008) Sterility Assurance Level (SAL): The term and its definition continues to cause confusion in the industry. *Pharm. Microbiol. Forum Newsl.* 14, 1–15.
- (99) Schleh, M., Romanowski, P., Bhebe, P., Zhang, L., Chinniah, S., Lawrence, B., Bashiri, H., Gaduh, A., Rajurs, V., Rasmussen, B., and Chuck, A. (2009) Susceptibility of Mouse Minute Virus to inactivation by heat in two cell culture media types. *Biotechnol. Prog.* 25, 854–860.
- (100) Weaver, B., and Rosenthal, S. (2010) Viral risk mitigation for mammalian cell culture media. *PDA J. Pharm. Sci. Technol.* 64, 436–439.
- (101) Pohlscheidt, M., Charaniya, S., Kulenovic, F., Corrales, M., Shiratori, M., Bourret, J., Meier, S., Fallon, E., and Kiss, R. (2014) Implementing high-temperature short-time media treatment in commercial-scale cell culture manufacturing processes. *Appl. Microbiol. Biotechnol.* 98, 2965–2971.
- (102) Cao, X., Stimpfl, G., Wen, Z.-Q., Frank, G., and Hunter, G. (2013) Identification and root cause analysis of cell culture media precipitates in the viral deactivation treatment with high-temperature/short-time method. *PDA J. Pharm. Sci. Technol.* 67, 63–73.
- (103) Hijnen, W., Beerendonk, E., and Medema, G. (2006) Inactivation credit of UV radiation for viruses, bacteria and protozoan (oo)cysts in water: a review. *Water Res.* 40, 3–22.
- (104) Abshire, R. (1988) Ultraviolet radiation: A method of sterilization in the pharmaceutical industry. *Ozone Sci. Eng.* 10, 25–38.
- (105) Schmidt, S., and Kauling, J. (2007) Process and laboratory scale UV inactivation of viruses and bacteria using an innovative coiled tube reactor. *Chem. Eng. Technol.* 30, 945–950.

- (106) Meng, Q., and Gerba, C. (1996) Comparative inactivation of enteric adenoviruses, poliovirus and coliphages by ultraviolet irradiation. *Water Res.* 30, 2665–2668.
- (107) Shanley, J. D. (1982) Ultraviolet irradiation of murine cytomegalovirus. *J. Gen. Virol.* 63, 251–4.
- (108) Lytle, C., and Sagripanti, J. (2005) Predicted inactivation of viruses of relevance to biodefense by solar radiation. *J. Virol.* 79, 14244–14252.
- (109) Harris, R., Coleman, P., and Morahan, P. (1974) Stability of minute virus of mice to chemical and physical agents. *Appl. Microbiol.* 28, 351–54.
- (110) Chevrefils, G., Caron, E., and Wright, H. (2006) UV dose required to achieve incremental log inactivation of bacteria, protozoa and viruses. *IUVA News* 8, 38–45.
- (111) Linden, K., Scheible, K., and Posy, P. (2011) Can UV protect the public from Adenovirus in drinking water? *Proc. Water Environ. Fed.* 8, 26–33.
- (112) Oms-Oliu, G., Martín-Belloso, O., and Soliva-Fortuny, R. (2008) Pulsed light treatments for food preservation. A review. *Food Bioprocess Technol.* 3, 13–23.
- (113) Ansari, I., and Datta, A. (2003) An overview of sterilization methods for packaging materials used in aseptic packaging systems. *Food Bioprod. Process.* 81, 57–65.
- (114) Bintsis, T., Litopoulour-Tzanetaki, E., and Robinson, R. K. (2000) Existing and potential applications of ultraviolet light in the food industry—a critical review. *J. Sci. Food Agric.* 80, 637–645.
- (115) Wekhof, A., Trompeter, F.-J., and Franken, O. (2001) Pulsed UV Disintegration (PUVD): A new sterilisation mechanism for packaging and broad medical-hospital applications, in *The First International Conference on Ultraviolet Technologies*, pp 1–15.
- (116) Yaman, A. (2001) Alternative methods of terminal sterilization for biologically active macromolecules. *Curr. Opin. Drug Discov. Devel.* 4, 760–3.
- (117) Agalloco, J., and Akers, J. (2010) Advanced Aseptic Processing Technology (Swarbrick, J., Ed.) 1st ed., pp 278–280. Informa Healthcare, London.
- (118) Rubbo, S. D., and Gardner, J. F. (1965) A Review of Sterilization and Disinfection, pp 1–96. Lloyd-Luke (Medical Books) Ltd., Chicago.

- (119) Lorenzen, A., James, C. A., and Kennedy, S. W. (1993) Effects of UV irradiation of cell culture medium on PCB-mediated porphyrin accumulation and EROD induction in chick embryo hepatocytes. *Toxicol. In Vitro* 7, 159–66.
- (120) Hart, R., and Boychyn, R. M. (2012) Cell Culture media for UVC exposure and methods related thereto. *US Pat. 20,120,214,204*. United States of America.
- (121) Qualls, R. G., and Johnson, J. D. (1983) Bioassay and dose measurement in UV disinfection. *Appl. Environ. Microbiol.* 45, 872–877.
- (122) Bolton, J. R., and Linden, K. G. (2003) Standardization of methods for fluence (UV dose) determination in bench-scale UV experiments. *J. Environ. Eng.* 129, 209–215.
- (123) Pirnie, M., Linden, K. G., and Malley, J. P. J. (2006) Ultraviolet disinfection guidance manual for the final long term 2 enhanced surface water treatment rule. *Environ. Prot.*, pp 1–436. Washington.
- (124) Morris, G. A., Barjat, H., and Horne, T. J. (1997) Reference deconvolution methods. *Prog. Nucl. Magn. Reson. Spectrosc.* 31, 197–256.
- (125) Adams, D., Korke, R., and Hu, W.-S. (2007) Application of stoichiometric and kinetic analysis to characterize cell growth and product formation, in *Animal Cell Biotechnology Methods and Protocols* (Portner, R., Ed.) 2nd ed., pp 269–284. Humana Press, Totowa, New Jersey.
- (126) Vásquez-Vivar, J., Denicola, A. D., Radi, R., and Augusto, O. (1997) Peroxynitrite-mediated decarboxylation of pyruvate to both carbon dioxide and carbon dioxide radical anion. *Chem. Res. Toxicol.* 10, 786–794.
- (127) Weil, L., Gordon, W. G., and Buchert, A. R. (1951) Photooxidation of amino acids in the presence of methylene blue. *Arch. Biochem. Biophys.* 33, 90–109.
- (128) Handbook, C. R. C. (2014) Vapor Pressure, in *CRC Handbook of Chemistry and Physics* (Haynes, W. M., Ed.) 95th ed., pp 88–117. CRC Press.
- (129) Jenner, A. M., Rafter, J., and Halliwell, B. (2005) Human fecal water content of phenolics: the extent of colonic exposure to aromatic compounds. *Free Radic. Biol. Med.* 38, 763–72.
- (130) Belenguer, A., Duncan, S. H., Holtrop, G., Anderson, S. E., Lobley, G. E., and Flint, H. J. (2007) Impact of pH on lactate formation and utilization by human fecal microbial communities. *Appl. Environ. Microbiol.* 73, 6526–33.

- (131) Walker, A. W., Duncan, S. H., McWilliam, E. C., Child, M. W., Flint, H. J., and Leitch, E. C. M. (2005) pH and Peptide Supply Can Radically Alter Bacterial Populations and Short-Chain Fatty Acid Ratios within Microbial Communities from the Human Colon pH and Peptide Supply Can Radically Alter Bacterial Populations and Short-Chain Fatty Acid Ratios within Mi.
- (132) Food and Drug Administration. (2001) Guidance for Industry, Bioanalytical Method Validation Guidance for Industry Bioanalytical Method Validation. Rockville, MD.
- (133) Elsdon, S. R., Hilton, M. G., and Waller, J. M. (1976) The end products of the metabolism of aromatic amino acids by Clostridia. *Arch. Microbiol.* 107, 283–8.
- (134) D'Ari, L., and Barker, H. (1985) p-Cresol formation by cell-free extracts of Clostridium difficile. *Arch. Microbiol.* 311–312.
- (135) Wishart, D. S., Tzur, D., Knox, C., Eisner, R., Guo, A. C., Young, N., Cheng, D., Jewell, K., Arndt, D., Sawhney, S., Fung, C., Nikolai, L., Lewis, M., Coutouly, M.-A., Forsythe, I., Tang, P., Shrivastava, S., Jeroncic, K., Stothard, P., Amegbey, G., Block, D., Hau, D. D., Wagner, J., Miniaci, J., Clements, M., Gebremedhin, M., Guo, N., Zhang, Y., Duggan, G. E., Macinnis, G. D., Weljie, A. M., Dowlatabadi, R., Bamforth, F., Clive, D., Greiner, R., Li, L., Marrie, T., Sykes, B. D., Vogel, H. J., and Querengesser, L. (2007) HMDB: the Human Metabolome Database. *Nucleic Acids Res.* 35, D521–6.
- (136) Bino, R. J., Hall, R. D., Fiehn, O., Kopka, J., Saito, K., Draper, J., Nikolau, B. J., Mendes, P., Roessner-Tunali, U., Beale, M. H., Trethewey, R. N., Lange, B. M., Wurtele, E. S., and Sumner, L. W. (2004) Potential of metabolomics as a functional genomics tool. *Trends Plant Sci.* 9, 418–25.
- (137) Tang, J. (2011) Microbial metabolomics. *Curr. Genomics* 12, 391–403.
- (138) Martins dos Santos, V. a P., and Damborsky, J. (2010) Systems biology at work. *Curr. Opin. Biotechnol.* 21, 498–501.
- (139) Claesson, M. J., Jeffery, I. B., Conde, S., Power, S. E., O'Connor, E. M., Cusack, S., Harris, H. M. B., Coakley, M., Lakshminarayanan, B., O'Sullivan, O., Fitzgerald, G. F., Deane, J., O'Connor, M., Harnedy, N., O'Connor, K., O'Mahony, D., van Sinderen, D., Wallace, M., Brennan, L., Stanton, C., Marchesi, J. R., Fitzgerald, A. P., Shanahan, F., Hill, C., Ross, R. P., and O'Toole, P. W. (2012) Gut microbiota composition correlates with diet and health in the elderly. *Nature* 488, 178–84.
- (140) Lenz, E. M., and Wilson, I. D. (2007) Analytical strategies in metabonomics. *J. Proteome Res.* 6, 443–58.

- (141) Dettmer, K. (2007) Mass spectrometry-based metabolomics. *Mass Spectrom. Rev.* 26, 51–78.
- (142) Ward, J. L., Baker, J. M., and Beale, M. H. (2007) Recent applications of NMR spectroscopy in plant metabolomics. *FEBS J.* 274, 1126–31.
- (143) Karunasena, E., McMahon, K. W., Chang, D., and Brashears, M. M. (2014) Host responses to the pathogen *Mycobacterium avium* subsp. *paratuberculosis* and beneficial microbes exhibit host sex specificity. *Appl. Environ. Microbiol.* 80, 4481–90.
- (144) Cinquin, C., Le Blay, G., Fliss, I., and Lacroix, C. (2004) Immobilization of infant fecal microbiota and utilization in an in vitro colonic fermentation model. *Microb. Ecol.* 48, 128–38.
- (145) Team, R. D. C. (2011) R: a language and environment for statistical computing. *R Found. Stat. Comput.*
- (146) Wickham, H. (2009) ggplot2: elegant graphics for data analysis (Gentleman, R., Hornik, K., and Parmigiani, G., Eds.), pp 1–202. Springer New York, New York, NY.
- (147) Petrof, E. O., Gloor, G. B., Vanner, S. J., Weese, S. J., Carter, D., Daigneault, M. C., Brown, E. M., Schroeter, K., and Allen-Vercoe, E. (2013) Stool substitute transplant therapy for the eradication of *Clostridium difficile* infection: “RePOOPulating” the gut. *Microbiome* 1, 1–12.
- (148) Buckel, W., and Barker, H. a. (1974) Two pathways of glutamate fermentation by anaerobic bacteria. *J. Bacteriol.* 117, 1248–60.
- (149) Barker, H. (1981) Amino acid degradation by anaerobic bacteria. *Annu. Rev. Biochem.* 50, 23–40.
- (150) Measures, J. (1975) Role of amino acids in osmoregulation of non-halophilic bacteria. *Nature* 257, 398–400.
- (151) Giaever, H., and Styrvold, O. (1988) Biochemical and genetic characterization of osmoregulatory trehalose synthesis in *Escherichia coli*. *J. Bacteriol.* 170, 2841–2849.
- (152) Yamamoto, I., Abe, A., Saito, H., and Ishimoto, M. (1984) The pathway of ammonia assimilation in *Bacteroides fragilis*. *J. Gen. Appl. Microbiol.* 30, 499–508.
- (153) Mortensen, P., Clausen, M., Bonnen, H., Hove, H., and Holtug, K. (1992) Colonic fermentation of ispaghula, wheat bran, glucose, and albumin to short-chain fatty acids and ammonia evaluated in vitro in 50 subjects. *J. Parenter. Enter. Nutr.* 16, 433–439.

- (154) Macfarlane, G. T., Gibson, G. R., and Cummings, J. H. (1992) Comparison of fermentation reactions in different regions of the human colon. *J. Appl. Bacteriol.* 72, 84–91.
- (155) White, D. (2007) The physiology and biochemistry of prokaryotes 3rd ed., pp 388–391. Oxford University Press Inc., New York, NY.
- (156) Macy, J., Ljungdahl, L., and Gottschalk, G. (1978) Pathway of succinate and propionate formation in *Bacteroides fragilis*. *J. Bacteriol.* 134, 84–91.
- (157) Burlingame, R., and Chapman, P. J. (1983) Catabolism of phenylpropionic acid and its 3-hydroxy derivative by *Escherichia coli*. *J. Bacteriol.* 155, 113–21.
- (158) Díaz, E., Ferrández, a, and García, J. L. (1998) Characterization of the hca cluster encoding the dioxygenolytic pathway for initial catabolism of 3-phenylpropionic acid in *Escherichia coli* K-12. *J. Bacteriol.* 180, 2915–23.
- (159) Moss, C., Lambert, M., and Goldsmith, D. (1970) Production of hydrocinnamic acid by clostridia. *Appl. Microbiol.* 19, 375–379.
- (160) Ploux, O., Soularue, P., Marquet, A., Gloeckler, R., and Lemoine, Y. (1992) Investigation of the first step of biotin biosynthesis in *Bacillus sphaericus*. Purification and characterization of the pimeloyl-CoA synthase, and uptake of pimelate. *Biochem. J.* 287, 685–90.
- (161) Rosner, A. (1975) Control of lysine biosynthesis in *Bacillus subtilis*: inhibition of diaminopimelate decarboxylase by lysine. *J. Bacteriol.* 121, 20–28.
- (162) Chapman, P., and Duggleby, R. (1967) Dicarboxylic acid catabolism by bacteria. *Biochem. J.* 103, 7–9.
- (163) Mart, N., Gibello, A., and Ndez, F. (1991) Catabolism of 3- and 4-hydroxyphenylacetic acid by *Klebsiella pneumoniae* 621–628.
- (164) Mohamed, M. E., Seyfried, B., Tschach, A., and Fuchs, G. (1993) Anaerobic oxidation of phenylacetate and 4-hydroxyphenylacetate to benzoyl-coenzyme A and CO₂ in denitrifying *Pseudomonas* sp. *Arch. Microbiol.* 159, 563–573.
- (165) Gibson, G. R., Cummings, J. H., and Macfarlane, G. T. (1991) Growth and activities of sulphate-reducing bacteria in gut contents of healthy subjects and patients with ulcerative colitis. *FEMS Microbiol. Lett.* 86, 103–112.

- (166) Smith, E., and Macfarlane, G. (1998) Enumeration of amino acid fermenting bacteria in the human large intestine: effects of pH and starch on peptide metabolism and dissimilation of amino acids. *FEMS Microbiol. Ecol.* 25, 355–368.
- (167) Hilton, M. G., Mead, G. C., and Elsdon, S. R. (1975) The metabolism of pyrimidines by proteolytic clostridia. *Arch. Microbiol.* 102, 145–9.
- (168) Goldfine, H., and Stadtman, E. R. (1960) Propionic acid metabolism: V. The conversion of β -alanine to propionic acid by cell-free extracts of *Clostridium propionicum*. *J. Biol. Chem.* 235, 2238–2245.
- (169) Hayaishi, O., Nishizuka, Y., Tatibana, M., Takeshita, M., and Kuno, S. (1961) Enzymatic studies on the metabolism of β -alanine. *J. Biol. Chem.* 236, 781–790.
- (170) Newmark, H. L., and Lupton, J. R. (1990) Determinants and consequences of colonic luminal pH : Implications for colon cancer. *Nutr. Cancer* 14, 161–173.
- (171) Schink, B. (1984) Fermentation of tartrate enantiomers by anaerobic bacteria, and description of two new species of strict anaerobes, *Ruminococcus pasteurii* and *Ilyobacter tartaricus*. *Arch. Microbiol.* 139, 409–414.
- (172) Chadwick, V., Vince, A., Killingley, M., and Wrong, O. (1978) The metabolism of tartrate in man and the rat. *Clin Sci Mol Med* 54, 273–281.
- (173) Hinton, A. J., and Hume, M. (1997) Inhibition of *Listeria monocytogenes* growth by *Veillonella* cultured on tartrate medium. *Clin. Infect. Dis.* 25, S120.
- (174) La Rivière, J. (1958) On the microbial metabolism of the tartaric acid isomers. Uitgeverij Waltman: Delft.
- (175) Deppenmeier, U., Hoffmeister, M., and Prust, C. (2002) Biochemistry and biotechnological applications of *Gluconobacter* strains. *Appl. Microbiol. Biotechnol.* 60, 233–42.
- (176) Chandrashekar, K., Felse, P., and Panda, T. (1999) Optimization of temperature and initial pH and kinetic analysis of tartaric acid production by *Gluconobacter suboxydans*. *Bioprocess Eng.* 20, 1–5.
- (177) Vogels, G. D., and Van der Drift, C. (1976) Degradation of purines and pyrimidines by microorganisms. *Bacteriol. Rev.* 40, 403–68.
- (178) Wolpert, E., Phillips, S., and Summerskill, W. (1971) Transport of urea and ammonia production in the human colon. *Lancet* 1387–1390.

- (179) Dai, Z.-L., Li, X.-L., Xi, P.-B., Zhang, J., Wu, G., and Zhu, W.-Y. (2013) L-Glutamine regulates amino acid utilization by intestinal bacteria. *Amino Acids* 45, 501–12.
- (180) Kristoffersen, T., and Nelson, F. E. (1955) Degradation of amino acids by *Lactobacillus casei* and some factors influencing deamination of serine. *Appl. Microbiol.* 3, 268–73.
- (181) Hayward, H. R., and Stadtman, T. C. (1959) Anaerobic degradation of choline I., *Vibrio cholerae* n. sp: Fermentation of Choline by an anaerobic, cytochrome-producing bacterium 78, 557–561.
- (182) Neill, A. R., Grime, D. W., and Dawson, R. M. (1978) Conversion of choline methyl groups through trimethylamine into methane in the rumen. *Biochem. J.* 170, 529–35.
- (183) Craciun, S., and Balskus, E. P. (2012) Microbial conversion of choline to trimethylamine requires a glycyl radical enzyme. *Proc. Natl. Acad. Sci. U. S. A.* 109, 21307–12.
- (184) Fink, R., McGaughey, C., Cline, R., and Fink, K. (1956) Metabolism of intermediate pyrimidine reduction products in vitro. *J. Biol. Chem.* 218, 1–7.
- (185) Lara, F. (1952) On the decomposition of pyrimidines by bacteria I.: Studies by means of the technique of simultaneous adaptation. *J. Bacteriol.* 64, 271–277.
- (186) Friedrich, M., and Schink, B. (1993) Hydrogen formation from glycolate driven by reversed electron transport in membrane vesicles of a syntrophic glycolate-oxidizing bacterium. *Eur. J. Biochem.* 217, 233–40.
- (187) Savage, G. M. (1974) Lincomycin and clindamycin: their role in chemotherapy of anaerobic and microaerophilic infections. *Postepy Hig. Med. Dosw.* 28, 573–85.
- (188) Chang, F., Sih, C., and Weisblum, B. (1966) Lincomycin, an inhibitor of aminoacyl sRNA binding to ribosomes. *Proc. Natl. Acad. Sci. U. S. A.* 55, 431–438.
- (189) Kasten, M. J. (1999) Clindamycin, metronidazole, and chloramphenicol. *Mayo Clin. Proc.* 74, 825–33.
- (190) Heimdahl, A., and Nord, C. (1982) Effect of erythromycin and clindamycin on the indigenous human anaerobic flora and new colonization of the gastrointestinal tract. *Eur. J. Clin. Microbiol.* 1, 38–48.

- (191) Centers for Disease Control and Prevention. (1995) Recommendations for preventing the spread of vancomycin resistance: recommendations of the Hospital Infection Control Practices Advisory Committee (HICPAC). *Morbidity Mortal. Wkly. Rep.* 44, 105–113.
- (192) Loll, P. J., Kaplan, J., Selinsky, B. S., and Axelsen, P. H. (1999) Vancomycin binding to low-affinity ligands: delineating a minimum set of interactions necessary for high-affinity binding. *J. Med. Chem.* 42, 4714–9.
- (193) Walsh, C. T. (1993) Vancomycin resistance: decoding the molecular logic. *Science* 261, 308–9.
- (194) Guyton, A. C., and Hall, J. E. I. (2006) Textbook of medical psychology 11th ed., pp 748–760. Elsevier Inc., Philadelphia.
- (195) Everly, G. S., and Lating, J. M. (2013) The anatomy and physiology of the human stress response, in *A clinical guide to the treatment of the human stress response* (Everly, G. S. J., and Lating, J. M., Eds.) 3rd ed., pp 15–48. Springer New York, New York, NY.
- (196) Eisenhofer, G., Aneman, A., Hooper, D., Rundqvist, B., and Friberg, P. (1996) Mesenteric organ production, hepatic metabolism, and renal elimination of norepinephrine and its metabolites in humans. *J. Neurochem.* 66, 1565–73.
- (197) Hawrelak, J. a, and Myers, S. P. (2004) The causes of intestinal dysbiosis: a review. *Altern. Med. Rev.* 9, 180–97.
- (198) Bailey, M., and Coe, C. (1999) Maternal separation disrupts the integrity of the intestinal microflora in infant rhesus monkeys. *Dev. Psychobiol.* 35, 146–155.
- (199) O'Donnell, P. M., Aviles, H., Lyte, M., and Sonnenfeld, G. (2006) Enhancement of in vitro growth of pathogenic bacteria by norepinephrine: importance of inoculum density and role of transferrin. *Appl. Environ. Microbiol.* 72, 5097–9.
- (200) Bailey, M. T. (2010) Psychological stress, immunity, and the effects on indigenous microflora, in *Microbial Endocrinology* (Lyte, M., and Freestone, P. P. E., Eds.), pp 191–212. Springer New York, New York, NY.
- (201) Lyte, M., and Ernst, S. (1992) Catecholamine induced growth of gram negative bacteria. *Life Sci.* 50, 203–212.
- (202) Yamamoto, I., Saito, H., and Ishimoto, M. (1987) Regulation of synthesis and reversible inactivation in vivo of dual coenzyme-specific glutamate dehydrogenase in *Bacteroides fragilis*. *J. Gen. Microbiol.* 133, 2773–80.

- (203) Gorbach, S. L., Spanknebel, G., Weinstein, L., Plaut, A. G., Nahas, L., and Levitan, R. (1969) Studies of intestinal microflora. VIII. Effect of lincomycin on the microbial population of the human intestine. *J. Infect. Dis.* 120, 298–304.
- (204) Edwards, C., Duerden, B., and Read, N. (1986) Effect of clindamycin on the ability of a continuous culture of colonic bacteria to ferment carbohydrate. *Gut* 27, 411–417.
- (205) Thauer, R. K., Jungermann, K., and Decker, K. (1977) Energy conservation in chemotrophic anaerobic bacteria. *Bacteriol. Rev.* 41, 809.
- (206) Lowe, S. E., Jain, M. K., and Zeikus, J. G. (1993) Biology, ecology, and biotechnological applications of anaerobic bacteria adapted to environmental stresses in temperature, pH, salinity, or substrates. *Microbiol. Rev.* 57, 451–509.
- (207) Hold, G. L., Schwiertz, A., Aminov, R., Blaut, M., and Flint, H. J. (2003) Oligonucleotide probes that detect quantitatively significant groups of butyrate-producing bacteria in human feces. *Appl. Environ. Microbiol.* 69, 4320–4324.
- (208) Yap, I. K., Li, J. V., Saric, J., Martin, F.-P., Davies, H., Yulan, W., D, W. I., Nicholson, J. K., Utzinger, J., Marchesi, J. R., and Holmes, E. (2008) Metabonomic and microbiological analysis of the dynamic effect of vancomycin-induced gut microbiota modification in the mouse. *J. Proteome Res.* 7, 3718–3728.
- (209) Woodmansey, E. (2004) Comparison of compositions and metabolic activities of fecal microbiotas in young adults and in antibiotic-treated and non-antibiotic-treated elderly subjects. *Appl. Environ. Microbiol.* 70, 6113–6122.
- (210) Zeikus, J. G. (1980) Chemical and fuel production by anaerobic bacteria. *Annu. Rev. Microbiol.* 34, 423–64.
- (211) Meiberg, J. B. M., and Harder, W. (1978) Aerobic and anaerobic metabolism of trimethylamine, dimethylamine and methylamine in *Hyphomicrobium* X. *J. Gen. Microbiol.* 106, 265–276.
- (212) Patterson, J., and Hespell, R. (1979) Trimethylamine and methylamine as growth substrates for rumen bacteria and *Methanosarcina barkeri*. *Curr. Microbiol.* 3, 79–83.
- (213) Ban, J., Vitale, L., and Kos, E. (1972) Thymine and uracil catabolism in *Escherichia coli*. *J. Gen. Microbiol.* 73, 267–72.
- (214) Gibbons, R. J., and Macdonald, J. B. (1960) Hemin and vitamin K compounds as required factors for the cultivation of certain strains of *Bacteroides melaninogenicus*. *J. Bacteriol.* 80, 164–70.

(215) Stojiljkovic, I., and Hantke, K. (1992) Hemin uptake system of *Yersinia enterocolitica*: similarities with other TonB-dependent systems in gram-negative bacteria. *EMBO J.* *11*, 4359–4367.

(216) Spízek, J., and Rezanka, T. (2004) Lincomycin, clindamycin and their applications. *Appl. Microbiol. Biotechnol.* *64*, 455–64.

Appendix A Supplementary Figures and Tables

A.1 Materials and Methods

Table A1. Cultured bacterial isolates for MET-1.⁶

Higher group	taxonomic	Closest species match	%	Relative abundance (by biomass)
Actinobacteria		Bifidobacterium adolescentis (two different strains)	99.79 99.79	1.5 1.5
		Bifidobacterium longum (two different strains)	99.86 99.16	0.5 2
		Collinsella aerofaciens	98.73	1
		Bacteroides ovatus	99.52	1.5
		Parabacteroides distasonis	99.45	1.5
Firmicutes	Bacilli	Lactobacillus casei/paracasei	99.47	1
		Lactobacillus casei	99.74	1
		Streptococcus mitis	99.79	0.75
	Clostridium cluster IV	Eubacterium desmolans	94.9	1
		Faecalibacterium prausnitzii	99.17	2
	Clostridium cluster VIII	Clostridium cocleatum	91.92	1
	Clostridium cluster IX	Acidaminococcus intestini	100	1
	Clostridium cluster XIVa	Blautia sp.	99.55	0.75
		Dorea longicatena	99.62	1
		(two different strains)	99.6	2
		Eubacterium eligens	98.15	2
		Eubacterium rectale (four different strains)	99.59 99.6 99.19 99.53	2 1.5 1 1
		Eubacterium ventriosum	100	1
		Lachnospira pectinoshiza	95.22	0.5
		Roseburia faecalis	99.65	1.5
		Roseburia intestinalis	100	1

⁶Adapted from Table 2.3 of Development of an In Vitro Fermentation Model to Culture the Human Distal Gut Microbiota⁹⁰. The species column refers to the closest species match, inferred by alignment of 16S rRNA sequence to GreenGenes database. The “%” column refers to the percent identity closest match. Units are presented as the equivalent of 1 X 10⁶ µL loopful of bacterial culture scraped from petri dish plates.

		Ruminococcus torques	99.15	1
		(two different strains)	99.29	2
		Ruminococcus sp. (two different strains)	97.35 98.4	1.5 0.5
	Clostridium cluster XV	Eubacterium limosum	97.05	1.5
Proteobacteria		Escherichia coli (two different strains)	99.8 99.6	0.5 -
		Raoultella sp.	99.4	1

Table A2. Cultured bacterial isolates for MET-2.⁷

Higher taxonomic group	Closest species match	%	Relative abundance (by biomass)
Actinobacteria	<i>Bifidobacterium longum</i>	100	1
	<i>Collinsella aerofaciens</i>	100	0.5
	<i>Microbacterium schleiferi</i>	99.34	1.5
	<i>Aldercreutzia equolifaciens</i>	99.76	1
	<i>Micrococcus luteus</i>	97.04	1
Bacteroidetes	<i>Bacteroides ovatus</i>	100	0.5
	<i>Parabacteroides merdae</i>	100	1
Firmicutes	<i>Eubacterium rectale</i>	100	1
	<i>Blautia luti</i>	98.91	1
	<i>Roseburia hominis</i>	99.04	1
	<i>Roseburia lactaris</i>	95.07	0.5
	<i>Ruminococcus albus</i>	96.96	1
	<i>Eubacterium eligens</i>	96.78	1
	<i>Ruminococcus torques</i> (two different strains)	99.27 100	1 1
	<i>Roseburia intestinalis</i>	100	1
	<i>Eubacterium fissicatena</i>	97.67	1
	<i>Eubacterium ventriosum</i>	97.37	1
	<i>Blautia coccooides</i>	99.85	0.25
	<i>Blautia hydrogenotrophica</i>	100	0.1
	<i>Dorea longicatena</i>	100	0.25
	<i>Dorea formicigenerans</i>	99.49	0.25
	<i>Clostridium ramosum</i>	96.14	0.25
	<i>Eubacterium limosum</i>	99.25	0.75
	<i>Streptococcus thermophilus</i>	100	0.25
	<i>Bacillus simplex</i>	98.7	0.25
	<i>Coprococcus catus</i>	99.19	0.25
	<i>Flavonifractor plautii</i>	96.21	0.25

⁷The species column refers to the closest species match, inferred by alignment of 16S rRNA sequence to GreenGenes database. The “%” column refers to the percent identity closest match. Units are presented as the equivalent of 1 X 10⁶ µL loopful of bacterial culture scraped from petri dish plates.

	<i>Streptococcus mitis</i>	100	0.3
	<i>Phascolarctobacterium sp.</i>	99.85	0.1
Proteobacteria	<i>Escherichia coli</i>	100	0.25
	<i>Parasutterella excrementihominis</i>	100	0.5

Table A3. Cultured bacterial isolates for MET-2A, MET-2B, MET-3A, and MET-3B.⁸

Higher taxonomic group	Closest species match	%	Relative abundance (by biomass)							
			Donor D				Donor A			
			Run 30		Run 33		Run 31		Run 32	
			MET -2A	MET -2B	MET -2A*	MET -2B*	MET -3A	MET -3B	MET -3A*	MET -3B*
Actinobacteria	<i>Bifidobacterium longum</i> (three different strains)	99.86 99.16 100.00	1 1.5 -	0.50 0.75 -	1 1.5 -	0.50 0.75 -	- - 1	- - 0.50	- - 1	- - 0.50
	<i>Bifidobacterium adolescentis</i> (two different strain) [†]	99.79 99.79	1 1	1 1	0.25 1	0.25 1	- -	- -	- -	- -
	<i>Collinsella aerofaciens</i> (two different strains)	98.73 100.00	0.50 -	1 -	0.50 -	1 -	- 0.50	- 1	- 0.50	- 1
	<i>Microbacterium schleiferi</i>	99.34	-	-	-	-	1.5	0.75	1.5	0.75
	<i>Aldercruzia equolfaciens</i> [†]	99.76	-	-	-	-	1	1	0.50	0.50
	<i>Micrococcus luteus</i>	97.04	-	-	-	-	1	1	1	1
Bacteroidetes	<i>Bacteroides ovatus</i> (two different strains)	99.52 100.00	0.50 -	1.5 -	0.50 -	1.5 -	- 0.50	- 1.5	- 0.50	- 1.5
	<i>Parabacteroides distasonis</i>	99.45	1	1	1	1	-	-	-	-
	<i>Parabacteroides merdae</i>	100.00	-	-	-	-	1	1	1	1
Firmicutes	<i>Eubacterium rectale</i> (five different strains)	99.59 99.60 99.19 99.53 100.00	1 1 1 0.50 -	0.25 0.25 1 1 -	1 1 1 0.50 -	0.25 0.25 1 1 -	- - - - 1	- - - - 0.25	- - - - 1	- - - - 0.25
	<i>Blautia luti</i>	98.91	-	-	-	-	1	0.25	1	0.25
	<i>Roseburia hominis</i>	99.04	-	-	-	-	1	1	1	1
	<i>Roseburia lactaris</i>	95.07	-	-	-	-	0.50	1	0.50	1
	<i>Faecalibacterium prausnitzii</i>	99.17	1	0.25	1	0.25	-	-	-	-

⁸Adapted from Table 2.3 of Adaptation of a Single-Stage Chemostat System to Model Functional Aspects of the Human Distal Gut Microbiota⁹¹. The species column refers to the closest species match, inferred by alignment of 16S rRNA sequence to GreenGenes database. The “%” column refers to the percent identity closest match. Units are presented as the equivalent of 1 X 10⁶ µL loopful of bacterial culture scraped from petri dish plates. [†] Bacterial isolates that differed between inocula from replicate runs. *Has slight deviation in bacterial composition compared to original defined communities in Run 30 and Run 31

	<i>Ruminococcus albus</i>	96.96	-	-	-	-	1	0.25	1	0.25
Firmicutes	<i>Eubacterium eligens</i> (2 different strains)	98.15 96.78	1 -	0.25 -	1 -	0.25 -	- 1	- 0.25	- 1	- 0.25
	<i>Ruminococcus torques</i> (four different strains)	99.15 99.29 99.27 100.00	1.5 1.5 - -	0.50 0.75 - -	1.5 1.5 - -	0.50 0.75 - -	- - 1.5 1.5	- - 0.50 0.75	- - 1.5 1.5	- - 0.50 0.75
	<i>Roseburia intestinalis</i> (two different strains)	100.00 100.00	1 -	1 -	1 -	1 -	- 1	- 1	- 1	- 1
	<i>Roseburia faecalis</i>	99.65	1	1	1	1	-	-	-	-
	<i>Eubacterium fissicatena</i>	97.67	-	-	-	-	1	1	1	1
	<i>Eubacterium ventriosum</i> (two different strains) [†]	100.00 97.37	1 -	1 -	0.10 -	0.10 -	- 1	- 1	- 0.25	- 0.25
	<i>Ruminococcus obeum</i> (two different strains)	99.55 97.35	0.50 0.25	1.5 1	0.50 0.25	1.5 1	- -	- -	- -	- -
	<i>Ruminococcus producta</i>	98.40	0.25	0.50	0.25	0.50	-	-	-	-
	<i>Blautia coccoides</i>	99.85	-	-	-	-	0.25	0.50	0.25	0.50
	<i>Blautia hydrogenotrophica</i>	100.00	-	-	-	-	0.10	0.50	0.10	0.50
	<i>Dorea longicatena</i> (three different strains)	99.62 99.60 100.00	0.25 0.25 -	1.5 1 -	0.25 0.25 -	1.5 1 -	- - 0.25	- - 1	- - 0.25	- - 1
	<i>Dorea formicigenerans</i>	99.49	-	-	-	-	0.25	1.5	0.25	1.5
	<i>Clostridium cocleatum</i> [†]	91.92	0.25	1	0.25	0.50	-	-	-	-
	<i>Clostridium ramosum</i>	96.14	-	-	-	-	0.25	1	0.25	1
	<i>Eubacterium limosum</i> (two different strains)	97.05 99.25	0.75 -	1.5 -	0.75 -	1.5 -	- 0.75	- 1.5	- 0.75	- 1.5
	<i>Lactobacillus casei</i>	99.74	0.25	1	0.25	1	-	-	-	-
	<i>Streptococcus thermophilus</i>	100.00	-	-	-	-	0.25	1	0.25	1
	<i>Lactobacillus paracasei</i>	99.47	0.25	1	0.25	1	-	-	-	-
	<i>Bacillus simplex</i>	98.70	-	-	-	-	0.25	1	0.25	1
Firmicutes	<i>Coprococcus catus</i>	99.19	-	-	-	-	0.25	1	0.25	1
	<i>Eubacterium</i>	94.90	0.25	1	0.25	1	-	-	-	-

	<i>desmolans</i>									
	<i>Flavonifractor plautii</i>	96.21	-	-	-	-	0.25	1	0.25	1
	<i>Streptococcus mitis</i> (two different strains)	99.79 100.00	0.30 -	0.30 -	0.30 -	0.30 -	- 0.30	- 0.30	- 0.30	- 0.30
	<i>Acidaminococcus intestini</i>	100.00	1.5	0.75	1.5	0.75	-	-	-	-
	<i>Phascolarctobacterium</i> sp.	99.85	-	-	-	-	0.10	0.05	0.10	0.05
Proteobacteria	<i>Escherichia coli</i> (two different strains)	99.80 100.00	0.25 -	0.50 -	0.25 -	0.50 -	- 0.25	- 0.50	- 0.25	- 0.50
	<i>Raoultella</i> sp.	99.40	0.50	0.25	0.50	0.25	-	-	-	-
	<i>Parasutterella excrementihominis</i>	100.00	-	-	-	-	0.5	0.25	0.5	0.25

A.2 Standardization of Fecal Water Samples and Analysis

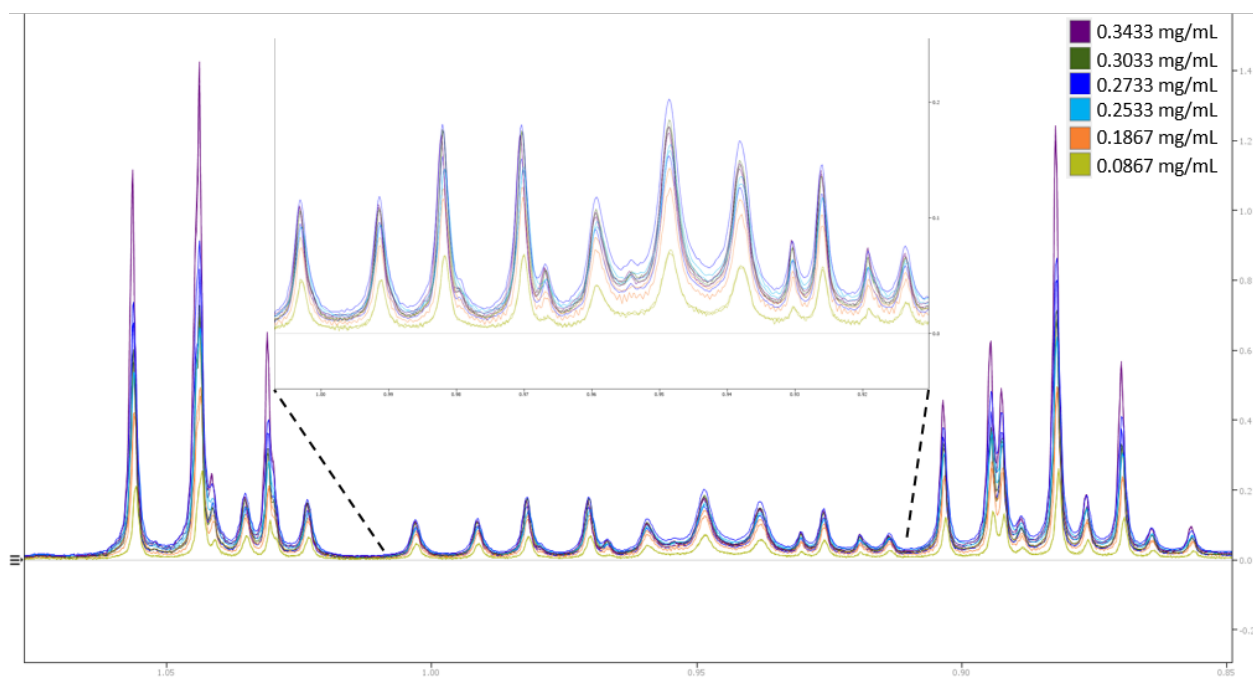


Figure A1. NMR Spectra of fecal water at various concentrations (0.85-1.05 ppm). Includes the NMR signature or partial signature of isobutyrate, propionate, isoleucine, valine, leucine, 2-oxoisocaproate, isovalerate, and butyrate from left to right. Inset is zoomed view of 0.90-1.01 ppm, includes NMR signatures or partial signatures of isoleucine, valine, leucine and 2-oxoisocaproate from left to right.

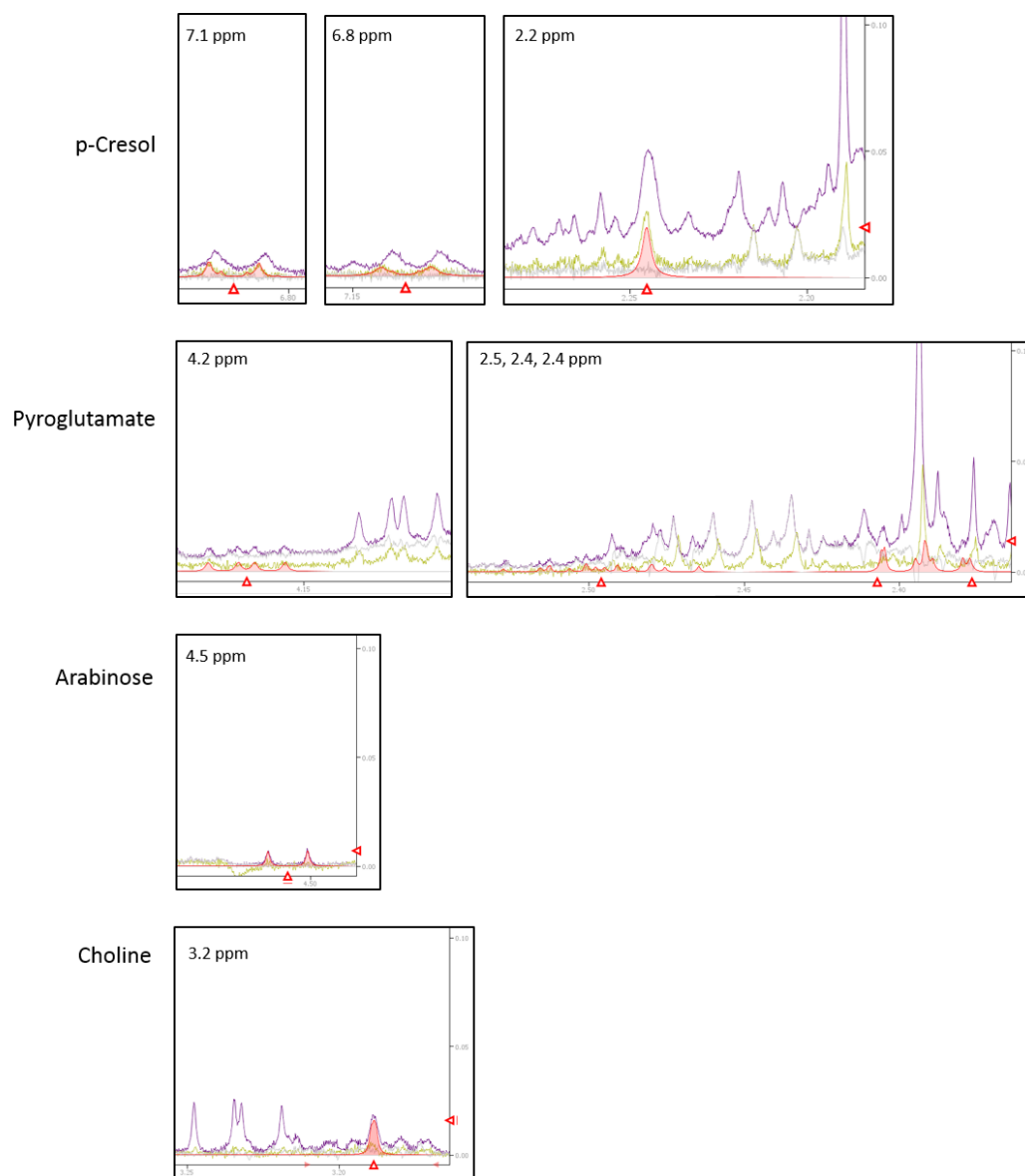


Figure A2. Profiling confidence impacted for certain compounds. Purple represents spectrum of 34.33% fecal water, and light green represents spectrum of 8.670% fecal water. For p-cresol, grey line represents the subtraction line between the 8.670% spectrum and area of profiled compounds in that spectrum. The p-Cresol signal, highlighted in red, represents the profiling of that compound fitted to the 8.670% sample spectrum. For all other compounds, the grey line represents the subtraction line between 34.33% fecal water and the area of all compounds profiled in that spectrum. The compound signals for pyroglutamate, arabinose and choline, highlighted in red, represent the profiling assignment of those compounds fitted to the 34.33% sample spectrum. Reason for demonstrating pyroglutamate, arabinose and choline relative to the high-concentration spectrum is because those compounds are not easily visualized in the low-concentration sample spectrum, due to low signal-to-noise ratio. Note: only compound clusters that were used in facilitating compound identification and quantifications are depicted.

A.3 Sample Preparation Procedure for Metabolomic Analysis

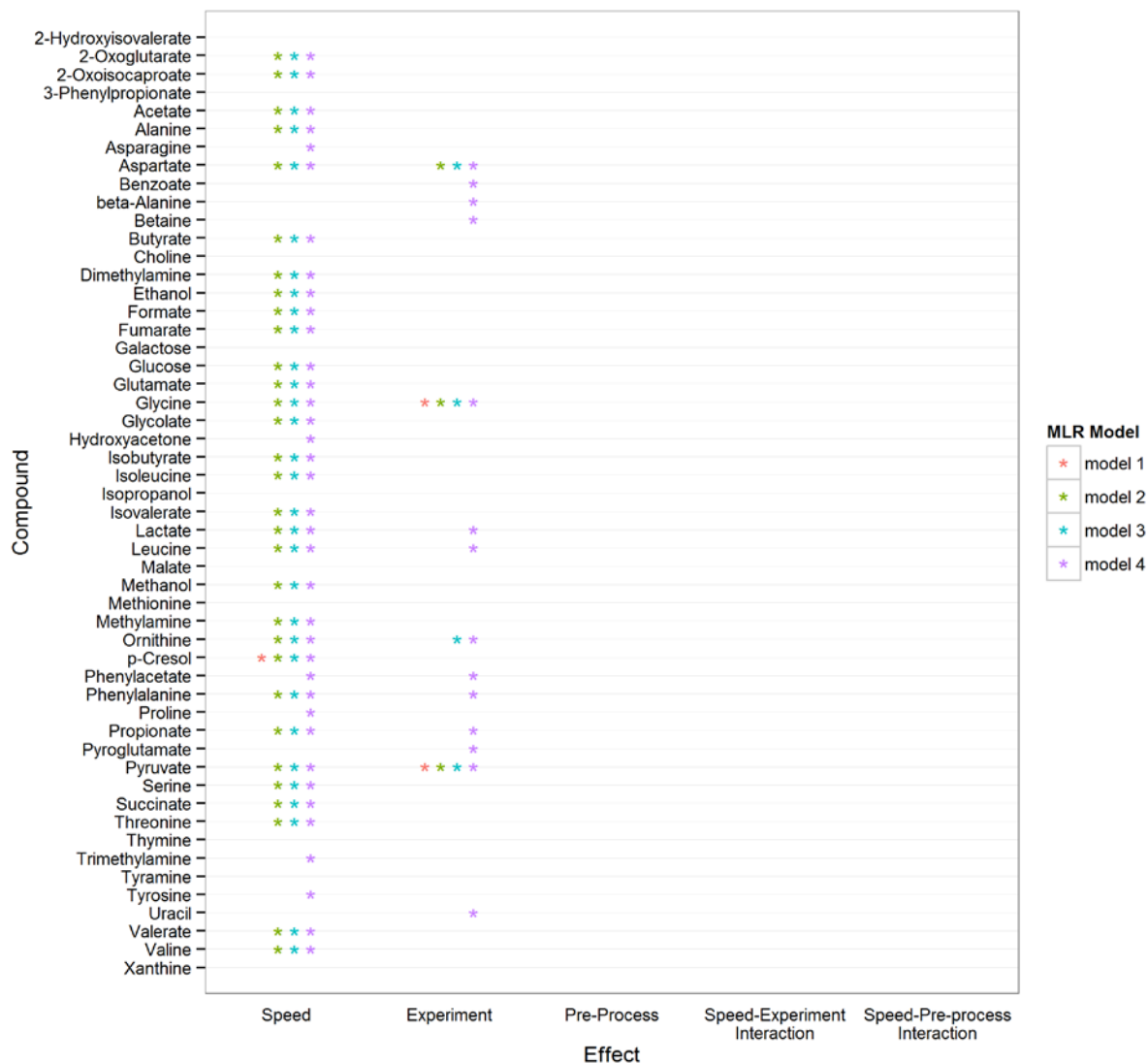


Figure A3. Significance of terms in each iteration of model simplification, with model 1 being the most complex model, and model 4 being the simplest model. Significant terms indicated with a star (*).

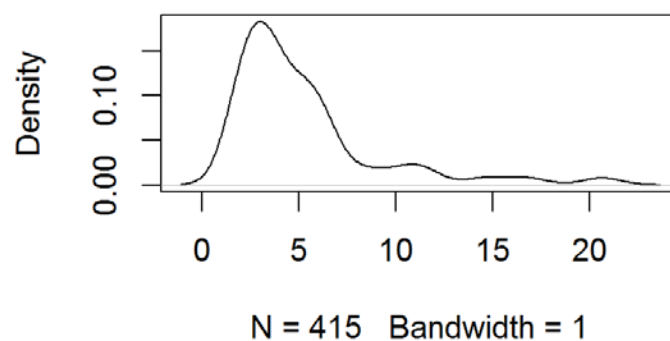


Figure A4. Evaluation of profiling precision based on relative standard deviation of compound concentrations observed in single-step filtered samples from Experiment 2.

Table A4. Compounds that did not demonstrate significant linear trend with changes in centrifugal force, based on simplest MLR model.

Compound	Slope	Standard error	T statistic	P value
3-Phenylpropionate	-9.61E-08	5.07E-08	-1.89346	0.064462
beta-Alanine	-4.98E-08	4.71E-08	-1.05815	0.295394
Choline	8.49E-10	3.59E-09	0.236667	0.813943
Pyroglutamate	-2.38E-07	1.48E-07	-1.61043	0.114001
Tyramine	2.39E-08	3.26E-08	0.734012	0.466588
Xanthine	-3.17E-08	3.69E-08	-0.85923	0.394576

Table A5. Compounds that demonstrated the same trends in concentration change with ultracentrifuge speed in both Experiment 1 and Experiment 2.

Compound	Experiment Bias	Standard Error	T Statistic	P-value
2-Hydroxyisovalerate	-0.00165	0.000561	-2.94227	0.005046
2-Oxoglutarate	-0.00253	0.001059	-2.38909	0.020956
2-Oxoisocaproate	-0.00024	0.000843	-0.28491	0.776966
3-Phenylpropionate	-0.00038	0.001903	-0.19831	0.843658
Acetate	2.170932	0.753433	2.881389	0.005949
Alanine	-0.01835	0.016419	-1.11733	0.269532
Asparagine	-0.00698	0.00222	-3.14535	0.002875
Butyrate	0.197922	0.089398	2.213936	0.031721
Choline	0.000274	0.000135	2.038151	0.047184
Dimethylamine	0.001538	0.000924	1.663707	0.102826
Ethanol	0.258782	0.102993	2.512618	0.015474
Formate	0.001515	0.002143	0.707101	0.483071
Fumarate	-6.57E-05	0.000249	-0.26356	0.793271
Galactose	-0.00855	0.004827	-1.77141	0.082977
Glucose	-0.00346	0.004161	-0.83057	0.410415
Glutamate	0.013564	0.014769	0.918433	0.363082
Glycolate	0.001112	0.002222	0.500279	0.619213
Hydroxyacetone	-0.00176	0.000749	-2.35505	0.022747
Isobutyrate	0.028285	0.016396	1.725171	0.091067
Isoleucine	-0.00475	0.010084	-0.47081	0.639952
Isopropanol	0.001274	0.002859	0.445497	0.658006
Isovalerate	-0.00559	0.008376	-0.66739	0.507788
Malate	0.002312	0.002725	0.848437	0.400495
Methanol	0.005318	0.002503	2.124866	0.038885
Methionine	0.001398	0.002021	0.691921	0.492391
Methylamine	0.006067	0.002394	2.53447	0.014652
p-Cresol	0.008506	0.003623	2.347627	0.023155
Proline	0.012527	0.003986	3.142632	0.002898
Serine	0.005399	0.003212	1.681038	0.09939
Succinate	0.065393	0.019404	3.37	0.00151
Threonine	6.22E-05	0.003425	0.018166	0.985584
Thymine	0.000269	0.00145	0.185337	0.853762
Trimethylamine	3.07E-05	0.000618	0.049737	0.960542
Tyramine	-0.00107	0.001222	-0.87233	0.387465
Tyrosine	0.013931	0.004123	3.378554	0.001473
Valerate	0.007871	0.028608	0.275138	0.784416

Valine	-0.02753	0.013235	-2.07976	0.043027
Xanthine	0.003677	0.001385	2.655878	0.010767

A.4 Effects of Varying Inoculant Biomass Proportions

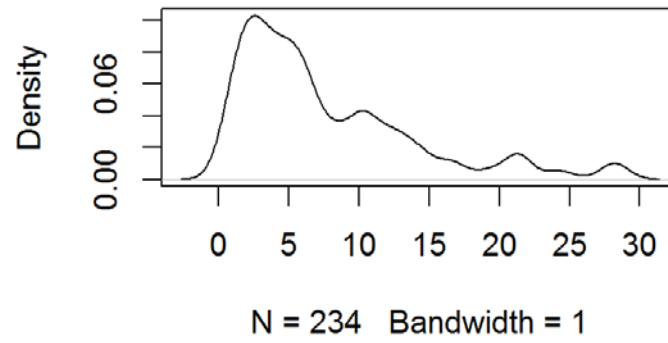


Figure A5. Density plot of the relative standard deviation of triplicate samples.

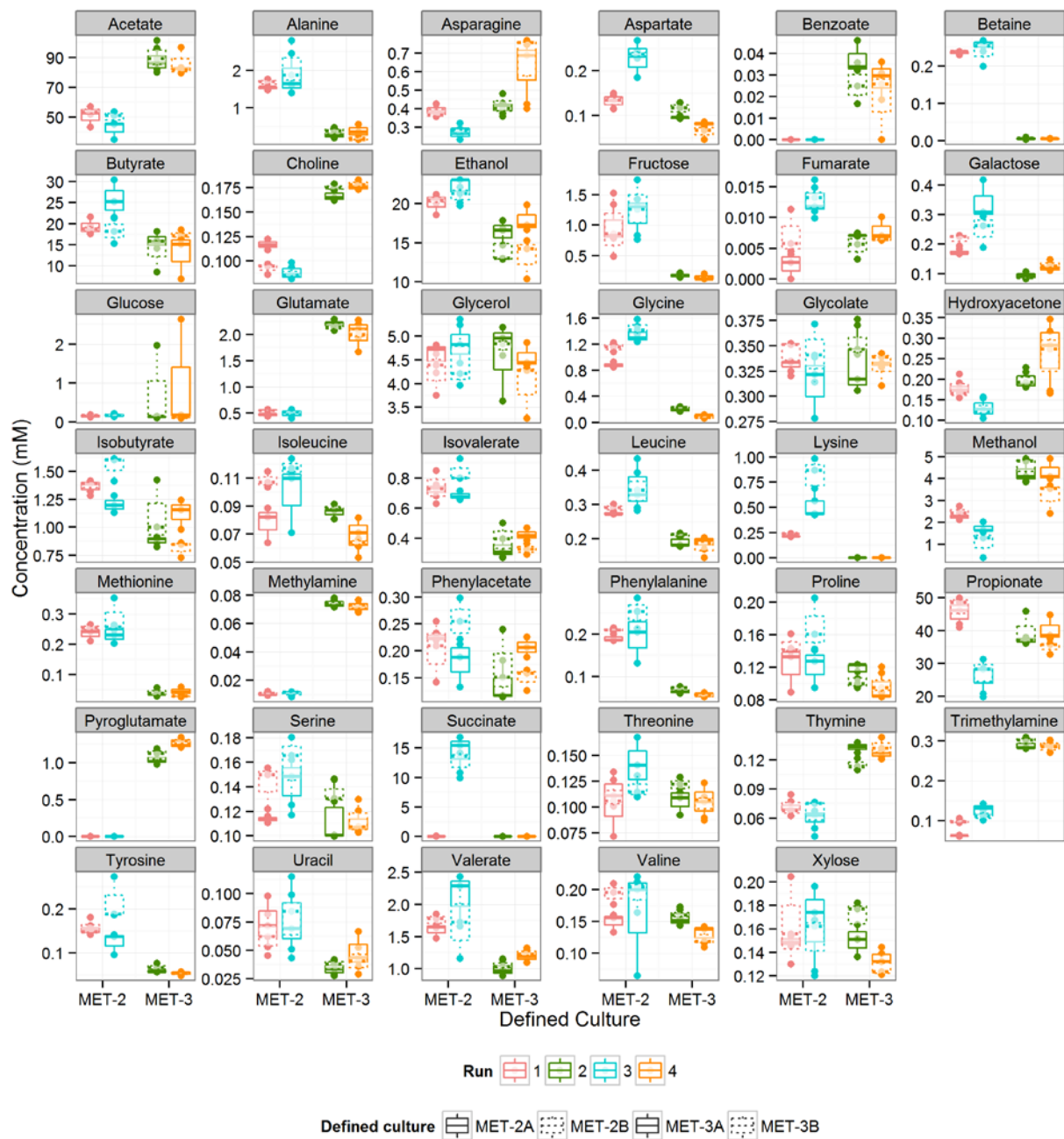


Figure A6. Profile overview of MET-2A and MET-3A and their reciprocal cultures, MET-2B and MET-3B, respectively.

A.5 METABOLOMIC ANALYSIS OF THE HUMAN GUT MICROBIOTA:
COMPARISON OF FECES-DERIVED COMMUNITIES AND DEFINED MIXED
COMMUNITIES

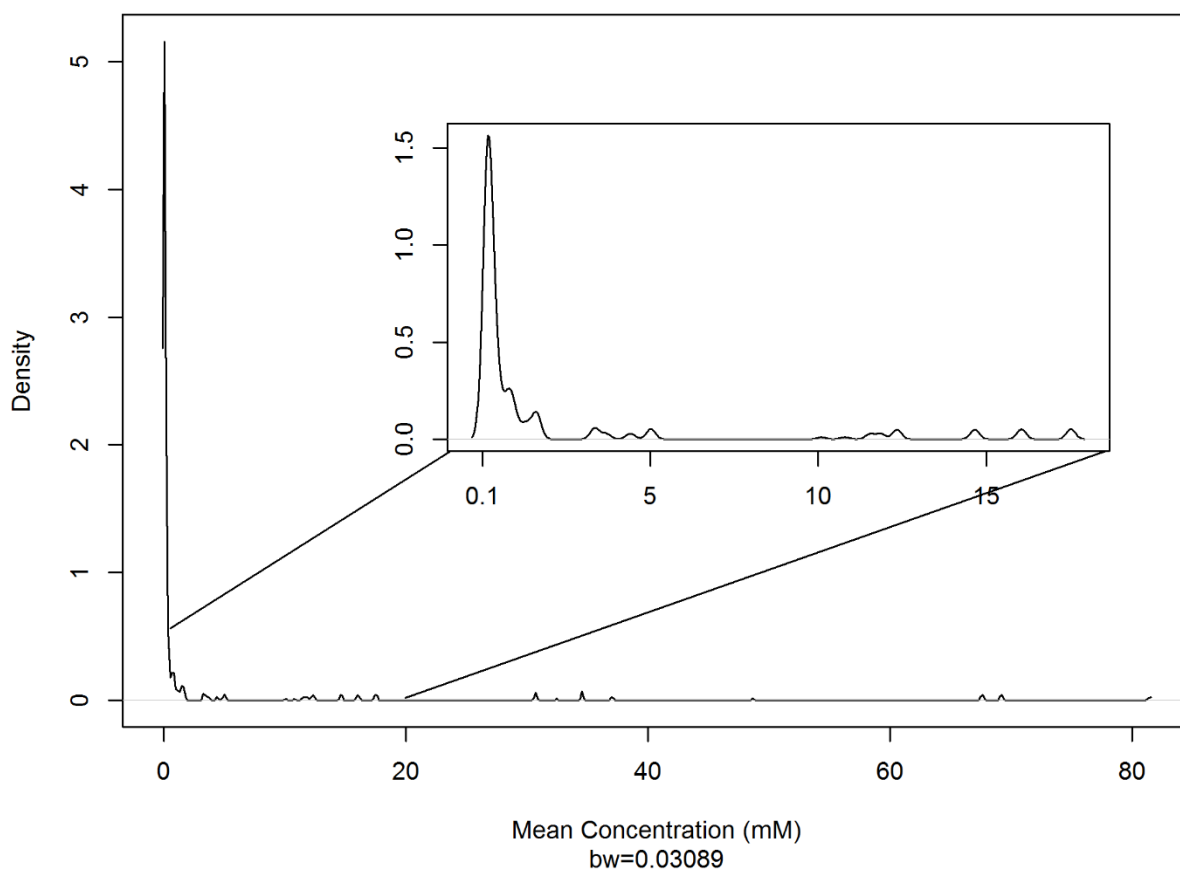


Figure A7. Kernel density plot of mean concentrations of untreated samples collected during steady state, bandwidth=0.03043. Inset is the kernel density plot of mean concentrations greater than or equal to 0.1 mM and less than or equal to 20 mM, bandwidth = 0.135. The clusters observed in these plots were used to define the bin in which the compounds would be grouped.

A.6 Analysis of Perturbations to Human Fecal Microbiota Propagated In Vitro

Table A6. List of microorganisms sensitive to clindamycin. Adapted from Spizek et al, 2004²¹⁶

Organism	Minimum Inhibitory concentration (µg/mL)
Gram-positive	
<i>Bacillus anthracis</i>	0.25–5.0
<i>Staphylococcus aureus</i>	0.04–1.6
<i>Sta. epidermidis</i>	0.1–0.2
<i>Streptococcus agalactiae</i>	0.02–0.1
<i>Str. pneumoniae</i>	0.002–0.1
<i>Str. pyogenes</i>	0.01–0.2
<i>Str. viridans</i>	0.005–0.2
Gram-negative	
<i>Escherichia coli</i>	64
<i>Haemophilus influenzae</i>	0.5–16.0
<i>Klebsiella pneumoniae</i>	125
<i>Neisseria gonorrhoeae</i>	0.5–4.0
<i>N. meningitis</i>	4
<i>Proteus vulgaris</i>	250
<i>Pseudomonas aeruginosa</i>	1000
<i>Salmonella schottmuelleri</i>	64

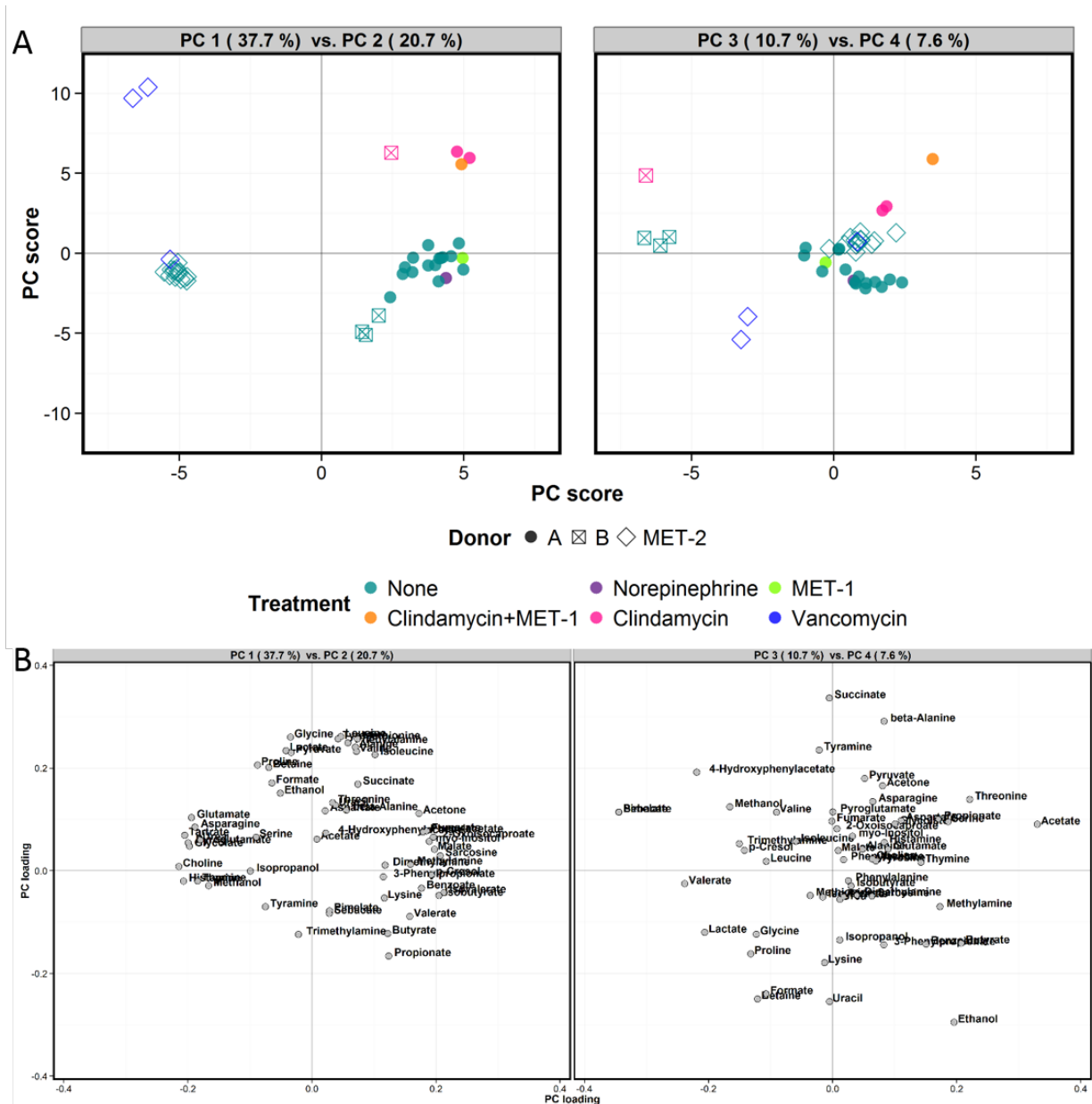


Figure A9. PCA model of metabolite profiles of Donor A, B, and MET-2, as well as profiles of each respective community culture after various treatments. (A) Score plot. (B) Loading plot.

Appendix B Calculations

Appendix B1. Relative Standard Deviation (RSD)

$$RSD = \frac{s}{\bar{x}} \times 100 \quad ; \quad s = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n - 1}}$$

Appendix B2. Equation for multiple linear regression.

The following equation was used to model each compound:

$$y_i = \beta_1 x_{i1} + \beta_2 x_2 + \beta_3 x_3 + \beta_3 x_1 x_2 + \beta_2 x_1 x_3 + \text{intercept} + \varepsilon$$

Where x_1 is ultracentrifuge speed, x_2 is experiment, x_3 is type of pre-processing step used, ε is the error and $i = 1 \dots n$ samples.

Appendix B3. Bonferroni correction on statistical significance when performing multiple comparisons.

$$\alpha_{corrected} = \frac{\alpha}{n}$$

Where n is the number of times the statistical test is applied.

Appendix B4. Percent change in concentration per 10 000 rpm

$$\text{percent change} = \left(\frac{\beta_1 * 10000}{y_{10000}} \right) \times 100$$

Where the response at rpm=10000, y_{10000} , and coefficient, β_1 , are from the MLR equation presented in Appendix B2.

Appendix B5. Analysis of Variance (ANOVA) about the mean. $MS_{residual}$ = Mean square within group, MS_{among} = Mean square between groups, n_j =number of observations within the j^{th} group, q =number of groups.

$$\bar{x}_{.j} = \frac{1}{n_j} \sum_{i=1}^{n_j} x_{ij} \quad ; \quad \bar{x}_{global} = \frac{1}{n} \sum_{j=1}^q \sum_{i=1}^{n_j} x_{ij} \quad ; \quad n = \sum_{j=1}^q n_j$$

$$SS_{residual} = \sum_i (x_{ij} - \bar{x}_{.j})^2 \quad ; \quad MS_{residual} = \frac{SS_{residual}}{n_j - q}$$

$$SS_{among} = \sum (\bar{x}_{.j} - \bar{x}_{global})^2 \quad ; \quad MS_{among} = \frac{SS_{among}}{q - 1}$$

$$F \text{ value} = \frac{MS_{residual}}{MS_{among}}$$

Appendix C Programming

C.1 Standardization of Fecal Water Samples and Analysis

Data source.R

```

1 #####
2 #Data source
3
4 require(RSQLite)
5 require(reshape2)
6
7 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
8 source('function-summary preprocessing.R')
9
10 setwd('C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R Scripts/Data source
    files')
11
12 #
13
14 raw <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R
    Scripts/Data source files/concentrations_x2-fecal water.csv')
15
16 d <- raw[['data']]
17 d.melt <- melt(d)
18 d.melt$value <- d.melt$value/0.9 #adjusting for sample dilution with 10% DSS
19
20 raw2 <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R
    Scripts/Data source files/concentrations_x1, x2, x3-fecal water.csv')
21
22 d2 <- raw2[['data']]
23 d2.melt <- melt(d2)
24 d2.melt$value <- d2.melt$value/0.9
25
26 raw3 <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R
    Scripts/Data source files/concentrations-fecal samples_Eric_SY.csv')
27
28 d3 <- raw3[['data']]
29 d3.melt <- melt(d3)
30 d3.melt$value <- d3.melt$value/0.9
31
32 #Qualifier
33 qualifier <- data.frame(file=d$file)
34 qualifier$sampleID <- gsub('.*([A-Z])(?=_x).*', '\\1', qualifier$sampleID, perl=
    TRUE)
35
36 qualifier$scan <- qualifier$file
37 qualifier$scan <- gsub('.*(?<=[A-Z]_x)([0-9]{1,2}).*', '\\1', qualifier$scan, perl=
    TRUE)
38
39 qualifier$profile.no <- qualifier$file
40 qualifier$profile.no <- gsub('.*([0-9])(?=\\.cnx).*', '\\1',
    qualifier$profile.no, perl=TRUE)
41
42
43 qualifier$sc.initial <- c(rep(8.2, 2), rep(5.6, 2), rep(9.1, 2), rep(20.6, 2),
    rep(15.2, 2), rep(2.6, 2))
44
45 qualifier$sv.initial <- c(rep(30,6), rep(60, 4), rep(30, 2))
46 qualifier$concentration <- round(qualifier$sc.initial/qualifier$sv.initial,4)
47
48 qualifier$rep <- c(rep(c(1,2),6))
49 qualifier$date <- 'June 5, 2014'
50
51 qualifier2 <- data.frame(file=d2$file)
52 qualifier2$sampleID <- qualifier2$file
53 qualifier2$sampleID <- gsub('.*([A-Z])(?=_x).*', '\\1', qualifier2$sampleID, perl=
    TRUE)
54

```

```

55 qualifier2$scan <- qualifier2$file
56 qualifier2$scan <- gsub('.*(?!=[A-Z]_x)([0-9]{1,2}).*', '\\1', qualifier2$scan,
57   perl=TRUE)
58 qualifier2$profile.no <- qualifier2$file
59 qualifier2$profile.no <- gsub('.*([0-9])(?=\\.cnx).*', '\\1',
60   qualifier2$profile.no, perl=TRUE)
61
62 qualifier2$sc.initial <- c(rep(c(rep(8.2, 2), rep(5.6, 2), rep(9.1, 2), rep(20.6, 2)
63   ,
64   rep(15.2, 2), rep(2.6, 2)),3))
65 qualifier2$v.initial <- c(rep(c(rep(30,6), rep(60, 4), rep(30, 2)),3))
66 qualifier2$concentration <- round(qualifier2$sc.initial/qualifier2$v.initial,4)
67 qualifier2$rep <- c(rep(c(rep(c(1,2),6)),3))
68 #Compound rank-----
69 rank <- read.csv('Compound rank-fecal water.csv')
70 rank <- rank[,1:4]
71 # Create the connection (no file needed) -----
72
73 # Choose driver
74 drv = dbDriver('SQLite')
75
76 # Create connection (for SQLite, this is a file)
77 con = dbConnect(drv, 'data.db')
78
79 # Creating database tables
80 data.table <- dbWriteTable(con, 'Data', d.melt, overwrite=TRUE)
81 data.table
82
83 data2.table <- dbWriteTable(con, 'Data2', d2.melt, overwrite=TRUE)
84 data2.table
85
86 data3.table <- dbWriteTable(con, 'Data3', d3.melt, overwrite=TRUE)
87 data3.table
88
89 qualifier.table <- dbWriteTable(con, 'Qualifier', qualifier, overwrite=TRUE)
90 qualifier.table
91
92 qualifier2.table <- dbWriteTable(con, 'Qualifier2', qualifier2, overwrite=TRUE)
93 qualifier2.table
94
95 rank.table <- dbWriteTable(con, 'Rank', rank, overwrite=TRUE)
96 rank.table
97
98 tables <- dbGetQuery(con, "SELECT * FROM sqlite_master WHERE type='table';")
99 print(tables)
100
101 dbDisconnect(con)

```

Fecal concentration analysis.R

```

1 #####
2 //
3 #Analysis
4 #loading required packages
5 require(ggplot2)
6 require(reshape2)
7 require(RSQLite)
8
9 #sourcing required functions
10 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
11 source('function-data cleanup.R')
12 source('function-descriptive statistics.R')
13
14 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation samples/Experiment 2/R
   Scripts/Functions')

```

```

15 source('correlation.R')
16
17 setwd('C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R Scripts/Analysis')
18
19 source('../Functions/linear model.R')
20
21 #-----
22 #Input data
23 # Choose driver
24 drv <- dbDriver('SQLite')
25
26 # Create connection (for SQLite, this is a file)
27 con <- dbConnect(drv, '../Data source files/data.db')
28
29 input <- dbGetQuery(con, "SELECT q.SampleID, q.concentration, q.rep,
30                          d.variable, d.value
31                          FROM Data d INNER JOIN Qualifier q
32                          ON d.file = q.file
33                          INNER JOIN Rank r
34                          ON d.variable = r.Compound
35                          WHERE r.Confidence = 1
36                          ORDER BY q.sampleID;")
37
38
39 dbDisconnect(con)
40
41 #Data clean up-----
42 clean <- data_cleanup(input, convert.NA=TRUE)
43 clean$rep <- as.character(clean$rep)
44
45 #Descriptive statistics-----
46 stat <- descriptive.stat(clean, data.type='long', variable=c('variable',
47                    concentration'))
48
49 #Overview-----
50 p <- ggplot(stat, aes(x=concentration, y=value)) +
51   geom_point(aes(colour=rep), alpha=0.8) +
52   xlab('Fecal water concentration (g/mL)') +
53   ylab('Compound concentration (mM)') +
54   facet_wrap(~variable, scales='free', ncol=6) +
55   theme_bw(14) +
56   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
57
58 ggsave('../Plots/Overview_x2.png', p, height=30, width=30, units='cm')
59
60 #Linear relationship-----
61 lm.test <- ddply(stat, 'variable', lm_custom, group='concentration', value='avg')
62 lm.sig <- lm.test[lm.test$slope.p <=0.05,]
63
64 #correlation
65 corr.result <- ddply(stat, 'variable', corr_custom, x='concentration', y='avg')
66 corr.sig <- corr.result[corr.result$p.val <= 0.05,]
67
68 plot.stat <- cbind(lm.test, corr.result$estimate)
69 colnames(plot.stat)[10] <- 'corr.estimate'
70 plot.stat$linetype <- 'ns'
71 plot.stat$linetype[plot.stat$slope.p <= 0.05 &
72   plot.stat$corr.estimate > 0.60] <- 's'
73
74 p.lm <- ggplot(stat, aes(x=concentration, y=value)) +
75   #geom_point(aes(colour=rep), alpha=0.6) +
76   geom_point(aes(y=avg), colour='black', alpha=0.8) +
77   geom_abline(data=plot.stat, aes(intercept=int, slope=slope, linetype=linetype)) +
78   scale_linetype_manual(values=c(2,1)) +
79   xlab('Fecal concentration (g/mL)') +
80   ylab('Compound concentration (mM)') +
81   facet_wrap(~variable, scales='free', ncol=6) +
82   theme_bw(9) +

```

```

82 theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
83
84 ggsave(' ../Plots/Relationship_x2_c.png', height=25, width=22, units='cm')
85
86 #Lowest concentrations
87 low <- stat[stat$concentration == 0.3433,]
88 low.cmpd <- unique(low[,c('variable','avg')])
89 low.cmpd <- low.cmpd[order(low.cmpd$avg),]

```

Normalization analysis.R

```

1 #####
2 ///
3 #normalizing for fecal water concentration
4
5 #loading required packages
6 require(ggplot2)
7 require(reshape2)
8 require(RSQLite)
9 require(gridExtra)
10
11 #sourcing required functions
12 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
13 source('function-data cleanup.R')
14 source('function-descriptive statistics.R')
15
16 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation samples/Experiment 2/R
17 Scripts/Functions')
18 source('correlation.R')
19
20 setwd('C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R Scripts/Analysis')
21
22 source('../Functions/linear model.R')
23 source('../Functions/stat custom.R')
24 source('../Functions/normalize.R')
25 #
26 #Input data
27 # Choose driver
28 drv <- dbDriver('SQLite')
29
30 # Create connection (for SQLite, this is a file)
31 con <- dbConnect(drv, '../Data source files/data.db')
32
33 input <- dbGetQuery(con, "SELECT q.SampleID, q.concentration, q.rep,
34 d.variable, d.value
35 FROM Data d INNER JOIN Qualifier q
36 ON d.file = q.file
37 INNER JOIN Rank r
38 ON d.variable = r.Compound
39 WHERE r.Confidence = 1
40 ORDER BY q.sampleID;")
41
42 input2 <- dbGetQuery(con, "SELECT * FROM Data3 d3 INNER JOIN Rank r
43 ON d3.variable = r.Compound
44 WHERE r.Confidence = 1;")
45
46 dbDisconnect(con)
47
48 #Data clean up
49 clean <- data_cleanup(input, convert.NA=TRUE)
50 clean$rep <- as.character(clean$rep)
51
52 #export in excel
53 cast <- dcast(clean, sampleID+concentration+rep~variable, value.var='value')
54 #write.csv(file='C:/Users/Sandi/Dropbox/Projects/Fecal water samples/R Scripts/
55 Results/fecal water data.csv', cast)
56
57 #Descriptive statistics

```

```

55 stat <- descriptive.stat(clean, data.type='long',
56                           variable=c('variable','concentration'))
57
58 #Linear relationship between metabolite concentration with fecal water
   concentration-----
59 lm.test <- ddply(stat, 'variable', lm_custom, group='concentration', value='avg')
60 lm.sig <- lm.test[lm.test$slope.p <=0.05,]
61
62 #correlation
63 corr.result <- ddply(stat, 'variable', corr_custom, x='concentration', y='avg')
64 corr.sig <- corr.result[corr.result$sp.val <= 0.05,]
65
66 plot.stat <- cbind(lm.test, corr.result$estimate)
67 colnames(plot.stat)[10] <- 'corr.estimate'
68 plot.stat$linetype <- 'ns'
69 plot.stat$linetype[plot.stat$slope.p <= 0.05 &
70                    plot.stat$corr.estimate > 0.60] <- 's'
71
72 p.lm <- ggplot(stat, aes(x=concentration, y=value)) +
73   geom_point(aes(colour=rep), alpha=0.6) +
74   geom_point(aes(y=avg), colour='black', alpha=0.8) +
75   geom_abline(data=plot.stat, aes(intercept=int, slope=slope, linetype=linetype)) +
76   scale_linetype_manual(values=c(2,1)) +
77   xlab('Fecal water concentration (g/mL)') +
78   ylab('Compound concentration (mM)') +
79   facet_wrap(~variable, scales='free', ncol=6) +
80   theme_bw(14) +
81   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
82
83 #Normalizing metabolite concentration for fecal water concentration-----
84 #expressing compound concentration as per g of feces
85 clean$norm <- clean$value/clean$concentration/1000
86
87 #Do the standardized concentrations have slope significantly different than zero
88 lm.norm <- ddply(clean, 'variable', lm_custom, group='concentration', value='norm')
89 lm.norm$linetype <- 'ns'
90 lm.norm$linetype[lm.norm$slope.p <=0.05] <- 's'
91
92 #Visualizing normalized values-----
93 p <- ggplot(clean, aes(x=concentration, y=norm)) +
94   geom_point(aes(colour=rep), alpha=0.8) +
95   geom_abline(data=lm.norm, aes(intercept=int, slope=slope, linetype=linetype)) +
96   scale_linetype_manual(values=c(2,1)) +
97   xlab('Fecal water concentration (g/mL)') +
98   ylab('Standardized Compound concentration (mmol/g feces)') +
99   facet_wrap(~variable, scales='free', ncol=6) +
100   theme_bw(14) +
101   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
102
103 #ggsave(' ../Plots/Normalized Overview.png', p, height=30, width=30, units='cm')
104
105
106 #looking at what the slopes of normalized values are
107 png(' ../Plots/Normalized slopes.png', units='cm', height=8, width=12, res=240)
108 plot(density(lm.norm$slope), main='', sub='N=39, Bandwidth=0.001',
109       xlab='Slope (mM/g of feces)')
110 dev.off()
111
112
113 p.norm <- ggplot(norm.stat, aes(x=variable, y=norm)) +
114   stat_boxplot() +
115   geom_point() +
116   xlab('Compound') +
117   ylab('Normalized Compound Concentration (mmol/g feces)') +
118   facet_wrap(~variable, scales='free', ncol=5) +
119   theme_bw(14) +
120   scale_x_discrete(labels=NULL) +
121   theme(axis.ticks.x=element_blank())

```

```

122 |
123 | ggsave(' ../Plots/Normalized concentration .b.png', p.norm,
124 |       height=35, width=30, units='cm')
125 |
126 | #Normalizing metabolite concentration given mM vs g/mL relationship as a score
127 | 

---


127 | clean$slope <- NA
128 | for(i in lm.test$variable) {
129 |   s <- lm.test$slope[lm.test$variable==i]
130 |   clean$slope[clean$variable == i] <- s
131 | }
132 |
133 | clean$score <- clean$value/clean$slope/clean$concentration/1000
134 |
135 | #Do the normalized scores have slope significantly different than zero
136 | lm.score <- ddply(clean, 'variable', lm_custom, group='concentration', value='score')
137 | lm.score$linetype <- 'ns'
138 | lm.score$linetype[lm.score$slope.p <=0.05] <- 's'
139 |
140 |
141 | p <- ggplot(clean, aes(x=concentration, y=score)) +
142 |   geom_point(aes(colour=rep), alpha=0.8) +
143 |   geom_abline(data=lm.score, aes(intercept=int, slope=slope, linetype=linetype)) +
144 |   scale_linetype_manual(values=c(2,1)) +
145 |   xlab('Fecal water concentration (g/mL)') +
146 |   ylab('Normalized Score of Compound Concentration') +
147 |   facet_wrap(~variable, scales='free', ncol=6) +
148 |   theme_bw(14) +
149 |   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
150 |
151 | #ggsave(' ../Plots/Normalized Score Overview.png', p, height=30, width=30, units='cm
152 | ')
153 | #Comparing Normalized profile to previous fecal sample profiles 

---


154 |
155 | clean2 <- data_cleanup(input2, convert.NA=TRUE)
156 | clean2$rep <- as.character(clean$rep)
157 | stat2 <- descriptive.stat(clean2, data.type='long', variable='variable')
158 |
159 | stat2$scaled <- stat2$value/100
160 | stat2$scaled[stat2$variable == 'Ethanol'] <- stat2$scaled[stat2$variable == '
161 |   Ethanol']/50
162 | stat2$expt <- 'Preliminary experiment'
163 | norm.stat$expt <- 'Smallest Poop experiment'
164 |
165 | pf <- ggplot(stat2, aes(x=variable, y=scaled, colour=expt), alpha=0.8) +
166 |   stat_boxplot(data=norm.stat, aes(y=norm)) +
167 |   geom_point(data=norm.stat, aes(y=norm)) +
168 |   stat_boxplot() +
169 |   geom_point() +
170 |   xlab('Compound') +
171 |   ylab('Compound Concentration') +
172 |   theme_bw(14) +
173 |   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
174 |         legend.position='none')
175 | #ggsave(' ../Plots/Old fecal samples.png', pf, height=15, width=35, units='cm')
176 | #Normalizing by highest concentration compound (Acetate) 

---


177 | cmpd.norm <- ddply(stat, 'concentration', norm_custom,
178 |                   norm='Acetate', norm.val='avg', norm.name='norm.acetate')
179 | cmpd.norm <- ddply(cmpd.norm, 'concentration', norm_custom,
180 |                   norm='Propionate', norm.val='avg', norm.name='norm.propionate')
181 | cmpd.norm <- ddply(cmpd.norm, 'concentration', norm_custom,
182 |                   norm='Butyrate', norm.val='avg', norm.name='norm.butyrate')
183 |
184 | norm <- melt(cmpd.norm, id.vars=c('sampleID', 'concentration', 'rep', 'variable'),
185 |             measure.vars=c('value', 'avg', 'norm.acetate', 'norm.propionate',
186 |                             'norm.butyrate'),

```

```

187     variable.name='norm.type')
188
189 norm <- norm[!norm$norm.type %in% c('value','avg'),]
190
191 cmpd.stat <- ddply(norm, 'norm.type', descriptive.stat, data.type='long',
192     variable=c('variable','concentration'))
193
194 #Compare feces norm and acetate normalization
195 #Getting descriptive stats of the normalized values
196 stat_func <- function(d, value){
197     stdev <- sd(d[,value])
198     avg <- mean(d[,value])
199     ymax <- max(d[,value])
200     ymin <- min(d[,value])
201     se <- stdev/sqrt(nrow(d))
202     out <- data.frame(avg=avg, stdev=stdev, se=se, ymax=ymax, ymin=ymin)
203     return(out)
204 }
205
206 fece.stat <- ddply(clean, 'variable', stat_func, value='norm')
207 acet.stat <- ddply(cmpd.norm, 'variable', stat_func, value='norm.acetate')
208
209 diff.stat <- data.frame(variable=fece.stat$variable, fec.sd=fece.stat$stdev,
210     fec.se=fece.stat$se, ace.sd=acet.stat$stdev,
211     ace.se=acet.stat$se)
212 diff.stat$diff.se <- diff.stat$fec.se-diff.stat$ace.se
213 diff.stat$diff.sd <- diff.stat$fec.sd-diff.stat$ace.sd
214
215 fece.stat$norm.type <- 'fec'
216 acet.stat$norm.type <- 'norm.acetate'
217
218 norm.stat <- rbind(fece.stat, acet.stat)
219
220 #Comparing Normalized profile to non-norm values
221 p.data <- clean
222 p.data$scal.val <- p.data$value/100
223 p.data <- melt(p.data, measure.vars=c('value','norm','scal.val'))
224 colnames(p.data)[5] <- 'data.type'
225 p.data <- p.data[p.data$data.type != 'value',]
226 #p.data <- merge(p.data, norm.stat[norm.stat$norm.type=='fec',], by='variable')
227
228 p.feces <- ggplot(p.data, aes(x=variable, y=value, colour=data.type)) +
229     stat_boxplot() +
230     #geom_point(position=position_dodge(width=1)) +
231     ylab('Compound Concentration') +
232     theme_bw(8) +
233     theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
234         axis.title.x=element_blank(),
235         legend.position=c(0.5,-0.35),
236         legend.direction='horizontal',
237         legend.margin=unit(0,'cm'),
238         panel.margin=unit(c(0.3,0.1,0,0.1),'cm'))
239
240 plot.data <- cmpd.stat[cmpd.stat$norm.type=='norm.acetate',]
241 p.acetate <- ggplot(plot.data, aes(x=variable, y=value)) +
242     stat_boxplot(aes(ymax=ymax, ymin=ymin)) +
243     #geom_point() +
244     theme_bw(8) +
245     xlab('Compound') +
246     ylab('Compound concentration (% Acetate)') +
247     theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
248         panel.margin=unit(c(0,0.1,0.2,0.1),'cm'))
249
250 g <- arrangeGrob(p.feces, p.acetate, ncol=1, nrow=2, height=c(0.5,0.5))
251
252 png(file='../Plots/norm_compare.png',width=18,height=18, units='cm',res=300)
253 print(g)
254 grid.text('A', x=0.0345,y=0.98)

```

```

255 grid.text('B', x=0.0345,y=0.51)
256 dev.off()
257
258 #Inner join two normalizations into one df for plotting
259 p.data <- data.frame(variable=clean$variable,
260                      concentration=clean$concentration,
261                      rep=clean$rep,
262                      fec=clean$norm)
263
264 p.data <- merge(p.data, compd.norm[,c('variable','concentration','rep',
265                                     'norm.acetate')],
266               by=c('variable','concentration','rep'))
267
268 p.data <- melt(p.data, id.vars=c('variable','concentration','rep'),
269               measure.vars=c('fec','norm.acetate'))
270
271 colnames(p.data)[4:5] <- c('norm.type','norm.val')
272
273 p.data <- merge(p.data, norm.stat, by=c('variable','norm.type'))
274
275 p.fec <- ggplot(p.data[p.data$norm.type=='fec',], aes(x=variable, y=norm.val)) +
276   geom_pointrange(aes(y=avg, ymin=ymin, ymax=ymax), size=0.3)+
277   #stat_boxplot() +
278   theme_bw(10) +
279   #xlab('Compound') +
280   ylab('Normalized concentration (mM/g feces)') +
281   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
282         axis.title.x=element_blank(),
283         plot.margin=unit(c(0.8,0.5,0,0.2),"cm"))
284
285 #ggsave('.. /Plots/norm range-feces.png', p.fec, height=11.5, width=18, units='cm')
286
287 p.ace <- ggplot(p.data[p.data$norm.type=='norm.acetate',], aes(x=variable, y=norm.
288   val)) +
289   geom_pointrange(aes(y=avg, ymin=ymin, ymax=ymax), size=0.3)+
290   #stat_boxplot() +
291   theme_bw(10) +
292   xlab('Compound') +
293   ylab('Normalized concentration (% acetate)') +
294   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
295         plot.margin=unit(c(0,0.5,0.2,0.2),"cm"))
296
297 #ggsave('.. /Plots/norm range-acetate.png', p.ace, height=11.5, width=18, units='cm')
298
299 g <- arrangeGrob(p.fec, p.ace, ncol=1, nrow=2, height=c(0.5,0.5))
300
301 png(file='../Plots/norm range.png',width=15,height=18, units='cm',res=300)
302 print(g)
303 grid.text('A', x=0.0345,y=0.98)
304 grid.text('B', x=0.0345,y=0.51)
305 dev.off()
306 #Compare difference compound normalizations
307
308 p.norm <- ggplot(compd.stat, aes(x=variable, y=value, colour=norm.type)) +
309   stat_boxplot(aes(ymax=ymax, ymin=ymin)) +
310   geom_point(position=dodge(width=0.8)) +
311   theme_bw(14) +
312   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
313
314 ggsave('.. /Plots/Normalize by diff compounds.png', p.norm,
315       height=15, width=35, units='cm')
316
317 #plotting only acetate normalized and poop normalized
318 plot.data <- compd.stat[compd.stat$norm.type=='norm.acetate',]
319 plot.data2 <- norm.stat
320 plot.data2$norm2 <- plot.data2$norm*25

```



```

321 p.acetate <- ggplot(plot.data, aes(x=variable, y=value)) +
322   #stat_boxplot(data=plot.data2, aes(y=norm2), colour='grey60') +
323   #geom_point(data=plot.data2, aes(y=norm2), colour='grey60') +
324   stat_boxplot(aes(ymax=ymax, ymin=ymin)) +
325   geom_point() +
326   theme_bw(14) +
327   xlab('Compound') +
328   ylab('Compound concentration (% Acetate)') +
329   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
330
331 ggsave('.. /Plots/Normalize by acetate.png', p.acetate,
332   height=15, width=35, units='cm')
333
334 pa2 <- ggplot(plot.data, aes(x=variable, y=value)) +
335   stat_boxplot(aes(ymax=ymax, ymin=ymin)) +
336   geom_point() +
337   facet_wrap(~variable, scales='free', ncol=5) +
338   theme_bw(14) +
339   xlab('Compound') +
340   ylab('Normalized Compound Concentration (% Acetate)') +
341   scale_x_discrete(labels=NULL) +
342   theme(axis.ticks.x=element_blank())
343
344 ggsave('.. /Plots/Normalize by acetate-concentrations2.png', pa2,
345   height=35, width=30, units='cm')

```

C.2 Gut-mimicking CSTR Effluent Sample Preparation and Analysisprep protocol

Data source.R

```

1 #####
2 #Data source for UC Expt 3
3
4
5 #sourcing required functions
6 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
7 source('function-summary preprocessing.R')
8
9 #loading required packages
10 require(reshape2)
11 require(RSQLite)
12
13 #Setting working directory
14 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
    Scripts/Data source files')
15
16 #preprocessing raw data-----
17 #Experiment 1
18 raw1 <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation
    Samples/Experiment 1/Results/R Scripts/Data source files/concentrations-UC Expt
    1.csv')
19 data1 <- raw1[['data']]
20 data.melt1 <- melt(data1)
21 data.melt1$value <- data.melt1$value/0.9
22
23 mw <- raw1[['mol.wt']]
24 mw1 <- as.data.frame(t(mw))
25 mw1$variable <- rownames(mw1)
26 mw1$Weight <- as.numeric(as.character(mw1$Weight))
27
28
29 #Experiment 2

```

```

30 raw2 <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation
    samples/Experiment 2/R Scripts/Data source files/concentrations_x2-UC Expt 2.
    csv')
31 data2 <- raw2[['data']]
32 data.melt2 <- melt(data2)
33 data.melt2$value <- data.melt2$value/0.9
34
35 mw <- raw2[['mol.wt']]
36 mw2 <- as.data.frame(t(mw))
37 mw2$variable <- rownames(mw2)
38 mw2$Weight <- as.numeric(as.character(mw2$Weight))
39
40 #Experiment 3
41 raw3 <- summary.preprocess('concentrations-UC Expt 3.csv')
42 data3 <- raw3[['data']]
43 data.melt3 <- melt(data3)
44 data.melt3$value <- data.melt3$value/0.9
45
46 mw <- raw1[['mol.wt']]
47 mw3 <- as.data.frame(t(mw))
48 mw3$variable <- rownames(mw3)
49 mw3$Weight <- as.numeric(as.character(mw3$Weight))
50
51 #Preliminary experiment
52 raw.prelim <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/
    Ultracentrifugation samples/Preliminary Experiment/R scripts/Data source files/
    D5-Ultracentrifugation concentration summary_2.csv')
53 data.prelim <- raw.prelim[['data']]
54 prelim.melt <- melt(data.prelim)
55 prelim.melt$value <- prelim.melt$value/0.9
56
57 mw <- raw1[['mol.wt']]
58 mw.prelim <- as.data.frame(t(mw))
59 mw.prelim$variable <- rownames(mw.prelim)
60 mw.prelim$Weight <- as.numeric(as.character(mw.prelim$Weight))
61
62 data <- rbind(data.melt1, data.melt2, data.melt3, prelim.melt)
63
64 mw <- unique(rbind(mw1, mw2, mw3, mw.prelim))
65 #Qualifiers-----
66 #file == file name from batch export
67 #scan == scan number within Expt
68 #sample == sample ID within UC Expt
69 #sampleID == sample ID within UC project -- unique across expts
70 #type == grouping triplicate samples within UC Expt -- not unique across expts
71 #grp == grouping triplicate samples across project (UC Expt 1 and 2 combined)
72 # -- unique across expts
73 #rep == replicate number
74 #class == processing procedure used
75 #speed == speed at which samples was ultracentrifuged
76 #process == presence of pre spin or no
77 #filter == filter size
78 #expt == experiment number
79
80 #Experiment 1-----
81 qualifier1 <- data.frame(file=data1$file)
82
83 qualifier1$scan <- qualifier1$file
84 qualifier1$scan <- gsub('.*(?!<=x)([0-9]{1,2}).*', '\\1',
85     qualifier1$scan, perl=TRUE)
86
87 scan <- read.csv('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/
    Experiment 1/Results/R Scripts/Data source files/scan order.csv', header=FALSE)
88 scan <- scan[2:nrow(scan),]
89 col.name <- as.character(as.matrix(scan[1,]))
90 colnames(scan) <- col.name
91 scan <- scan[2:nrow(scan),]
92 scan[, 'Scan Order'] <- as.numeric(scan[, 'Scan Order'])

```

```

93 scan$sampleID <- as.character(scan$sample)
94
95 type <- c()
96 for(i in 1:9) {
97   temp <- rep(LETTERS[i], 3)
98   type <- c(type, temp)
99 }
100
101 grp <- c()
102 for(i in 1:18) {
103   temp <- rep(LETTERS[i], 3)
104   grp <- c(grp, temp)
105 }
106
107 scan$type <- type
108 scan$grp <- grp[1:nrow(scan)]
109 scan$rep <- c(rep(1:3, 9))
110 scan$class <- c(rep('I',15), rep('II',6), rep('III',6))
111 scan$speed <- c(rep(10000,3), rep(20000,3), rep(30000,3), rep(40000,3), rep
    (50000,3),
112   rep(10000,3), rep(50000,3),rep(0,6))
113 scan$process <- c(rep('pre-spin',15), rep('no pre-spin',6), rep('pre-spin',6))
114 scan$filter <- c(rep('0.22',21), rep('serial',3),rep('0.22',3))
115
116 scan <- scan[order(scan[, 'Scan Order']),]
117
118 qualifier1 <- cbind(qualifier1, scan)
119 qualifier1 <- qualifier1[, !colnames(qualifier1) %in% c('sample_name',
120   'Scan Order')]
121
122 qualifier1$expt <- 'expt1'
123 qualifier1$sampleID <- factor(qualifier1$sampleID,
124   levels=c(LETTERS[1:26],
125   paste(LETTERS[1],LETTERS[1],sep='')))
126 qualifier1 <- qualifier1[order(qualifier1$sampleID),]
127
128 #Experiment 2
129 qualifier2 <- data.frame(file=data2$file)
130 qualifier2$scan <- qualifier2$file
131 qualifier2$scan <- gsub('.*([0-9]{2})(?=_x2\\.cnx).*', '\\1', qualifier2$scan, perl
    =TRUE)
132 qualifier2$scan <- gsub('.*([0-9]{1})(?=_x2\\.cnx).*', '\\1', qualifier2$scan, perl
    =TRUE)
133 qualifier2$sample <- qualifier2$file
134 qualifier2$sample <- gsub('([A-Z]{1,2})(?=_x[0-9]{1,2}_).*', '\\1',
135   qualifier2$sample, perl=TRUE)
136
137 qualifier2$sample <- factor(qualifier2$sample,
138   levels=c(LETTERS[1:26],
139   paste(LETTERS[1],LETTERS[1],sep='')))
140 qualifier2 <- qualifier2[order(qualifier2$sample),]
141
142
143 sampleID <- c(paste(LETTERS[1],LETTERS[2:26], sep=''),
144   paste(LETTERS[2],LETTERS[1:26], sep=''),
145   paste(LETTERS[3],LETTERS[1:26], sep=''))
146
147 qualifier2$sampleID <- sampleID[1:nrow(qualifier2)]
148
149 qualifier2$type <- type
150 qualifier2$grp <- grp[nrow(scan)+1:nrow(qualifier2)]
151 qualifier2$rep <- c(rep(1:3, 9))
152 qualifier2$class <- c(rep('I',15), rep('II',6), rep('III',6))
153 qualifier2$speed <- c(rep(10000,3), rep(20000,3), rep(30000,3), rep(40000,3), rep
    (50000,3),
154   rep(10000,3), rep(50000,3),rep(0,6))
155 qualifier2$process <- c(rep('pre-spin',15), rep('no pre-spin',6), rep('pre-spin',
    6))

```

```

156 qualifier2$filter <- c(rep('0.22',21), rep('serial',3),rep('0.22',3))
157 qualifier2$expt <- 'expt2'
158
159 #Experiment 3
160 qualifier3 <- data.frame(file=data3$file, scan=data3$file)
161 qualifier3$scan <- gsub('.*(?!<=x)([0-9]{1,2}).*', '\\1', qualifier3$scan,
162 perl=TRUE)
163
164 scan.order <- read.csv('scan_order.csv')
165 scan.order <- scan.order[,2:3]
166
167 qualifier3$sample <- scan.order[,2]
168 qualifier3 <- qualifier3[order(qualifier3$sample),]
169
170 qualifier3$sampleID <- sampleID[(nrow(qualifier2)+1):
171 (nrow(qualifier2)+nrow(qualifier3))]
172
173 qualifier3$type <- c(rep('A',8), unlist(lapply(LETTERS[2:6], rep, 3)))
174 qualifier3$grp <- c(rep('S',8), unlist(lapply(LETTERS[20:24], rep, 3)))
175 qualifier3$rep <- c(1:8, rep(1:3,5))
176 qualifier3$class <- c(rep('III',8), rep('I',15))
177
178 sample.info <- lapply(seq(50000, 10000, -10000), rep, times=3)
179 qualifier3$speed <- c(rep(0,8), unlist(sample.info))
180
181 qualifier3$process <- 'no pre-spin'
182 qualifier3$filter <- c(rep('serial',8), rep(0.22,15))
183 qualifier3$expt <- 'expt3'
184
185 #Preliminary experiment
186 qualifier.prelim <- data.frame(file=data.prelim$file)
187 qualifier.prelim$scan <- c(1,2)
188 qualifier.prelim$sample <- LETTERS[1:nrow(qualifier.prelim)]
189 qualifier.prelim$sampleID <- sampleID[(nrow(qualifier2)+nrow(qualifier3)+1):
190 (nrow(qualifier2)+nrow(qualifier3)+
191 nrow(qualifier.prelim))]
192 qualifier.prelim$type <- LETTERS[1:length(unique(data.prelim$file))]
193 qualifier.prelim$grp <- LETTERS[(length(unique(qualifier1$grp))+
194 length(unique(qualifier2$grp))+1):
195 (length(unique(qualifier1$grp))+
196 length(unique(qualifier2$grp))+
197 length(unique(qualifier.prelim$file)))]
198 qualifier.prelim$rep <- rep(1,2)
199 qualifier.prelim$class <- c('II', 'III')
200 qualifier.prelim$speed <- c(50000, 0)
201 qualifier.prelim$process <- c('no pre-spin', 'pre-spin')
202 qualifier.prelim$filter <- c(0.22, 'serial')
203 qualifier.prelim$expt <- 'preliminary'
204
205 #
206 qualifier <- rbind(qualifier1, qualifier2, qualifier3, qualifier.prelim)
207 qualifier$sampleID <- as.character(qualifier$sampleID)
208
209 #Rank
210 #Experiment 1
211 rank1 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/
212 Experiment 1/Results/R Scripts/Data source files/compound ranking.csv',
213 encoding='UTF-8')
214 rank1 <- rank1[,1:4]
215 rank1$expt <- 'expt1'
216
217 #Experiment 2
218 rank2 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation samples/
219 Experiment 2/R Scripts/Data source files/Compound rank-UC Expt 2.csv', encoding
220 ='UTF-8')
221 rank2 <- rank2[,1:4]
222 rank2$expt <- 'expt2'
223

```

```

220 #Experiment 3
221 rank3 <- read.csv('Compound rank-UC Expt 3.csv')
222 rank3 <- rank3[,1:4]
223 rank3$expt <- 'expt3'
224
225 #Preliminary experiment
226 rank.prelim <- read.csv('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation
      samples/Preliminary Experiment/R scripts/Data source files/Compound annotation.
      csv', encoding='UTF-8')
227 rank.prelim <- rank.prelim[,c(1:3,5)]
228 rank.prelim$expt <- 'preliminary'
229
230 rank <- rbind(rank1, rank2, rank3, rank.prelim)
231
232 #Compound info-----
233 info <- read.csv('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Data
      source files/Database-Compound classes.csv')
234
235 # Create the connection (no file needed) -----
236
237 # Choose driver
238 drv = dbDriver('SQLite')
239
240 # Create connection (for SQLite, this is a file)
241 con = dbConnect(drv, 'data.db')
242
243 # Creating database tables
244 uc.table <- dbWriteTable(con, 'UC3_data', data, overwrite=TRUE)
245 uc.table
246
247 qualifier.table <- dbWriteTable(con, 'Qualifier', qualifier, overwrite=TRUE)
248 qualifier.table
249
250 rank.table <- dbWriteTable(con, 'Rank', rank, overwrite=TRUE)
251 rank.table
252
253 info.table <- dbWriteTable(con, 'Cmpd_info', info, overwrite=TRUE)
254 info.table
255
256 mw.table <- dbWriteTable(con, 'Mol_wt', mw, overwrite=TRUE)
257 mw.table
258
259 tables <- dbGetQuery(con, "SELECT * FROM sqlite_master WHERE type='table';")
260 print(tables)
261
262 dbDisconnect(con)

```

Profiling precision analysis.R

```

1 #####
2 #Ultracentrifuge Experiment 3 - Analysis of profiling variability
3
4 #Pre-ambles-----
5 #Setting working directory
6 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Data source
      files/')
7
8 #Sourcing required functions
9 source('..../Functions/function-data cleanup.R')
10 source('..../Functions/function-descriptive statistics.R')
11
12 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 2/R
      Scripts/Functions')
13 source('..../ANOVA-mean.R')
14 source('..../kruskal wallis.R')
15 source('..../linear model.R')
16 source('..../tukey.R')
17 source('..../correlation.R')

```

```

18
19 #loading required packages
20 require(RSQLite)
21 require(reshape2)
22 require(plyr)
23 require(ggplot2)
24
25 #setting working directory
26 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
    Scripts/Analysis')
27 #
28 #Input data
29 # Choose driver
30 drv <- dbDriver('SQLite')
31
32 # Create connection (for SQLite, this is a file)
33 con <- dbConnect(drv, '../Data source files/data.db')
34
35 expt3 <- dbGetQuery(con, "SELECT q.sampleID, q.type, q.grp, q.rep, q.class,
36                          q.speed, q.process, q.filter, q.expt, c.Type, m.Weight,
37                          d.variable, d.value
38                          FROM UC3_Data d INNER JOIN Qualifier q
39                          ON d.file = q.file
40                          INNER JOIN Rank r
41                          ON d.variable = r.Compound AND q.expt = r.expt
42                          INNER JOIN Cmpd_info c
43                          ON d.variable = c.Compound
44                          INNER JOIN Mol_wt m
45                          ON d.variable = m.variable
46                          WHERE r.Confidence = 1 AND q.expt = 'expt3'
47                          AND q.type = 'A'
48                          AND NOT (d.variable = 'Formate' AND d.value > 0.3)
49                          ORDER BY q.type;")
50
51 dbDisconnect(con)
52
53 #Data clean up
54 clean <- data_cleanup(expt3, convert.NA=TRUE)
55
56 #Descriptive statistics
57 stat <- descriptive.stat(clean, data.type='long', variable=c('type', 'variable'))
58
59 #Relative standard deviation
60
61 stat$rsd <- stat$stdev/stat$avg*100
62
63 plot(density(stat$rsd, bw=1))
64
65 png(filename='../Plots/precision-RSD.png',
66      width=10, height=7, units='cm', res=240)
67 plot(density(stat$rsd, bw=1), main='')
68 dev.off()
69
70 high.rsd <- unique(stat[stat$rsd > 15, c('variable', 'rsd')])
71 write.csv(file='../Results/precision-rsd.csv', high.rsd)

```

Factor effects-mlr analysis.R

```

1 #####
2 //
3 #Comparison of profiles from Expt 3, and Expt 2
4 #looking at effects of all explanatory variables (process, speed, expt)
5
6 #pre-amble
7 #Setting working directory
8 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
9
10 #Sourcing required functions

```

```

10 source('function-data cleanup.R')
11 source('function-descriptive statistics.R')
12
13 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 2/R
    Scripts/Functions')
14 source('./ANOVA-mean.R')
15 source('./kruskal wallis.R')
16 source('./linear model.R')
17 source('./tukey.R')
18 source('./correlation.R')
19
20 #loading required packages
21 require(RSQLite)
22 require(reshape2)
23 require(plyr)
24 require(ggplot2)
25
26 #setting working directory
27 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
    Scripts/Analysis')
28 source('../Functions/multiple linear regression.R')
29 source('../Functions/mlr prediction.R')
30 #

```

```

31 #Input data
32 # Choose driver
33 drv <- dbDriver('SQLite')
34
35 # Create connection (for SQLite, this is a file)
36 con <- dbConnect(drv, '../Data source files/data.db')
37
38 input <- dbGetQuery(con, "SELECT q.SampleID, q.type, q.grp, q.rep, q.class,
39                          q.speed, q.process, q.filter, q.expt, c.Type, m.Weight,
40                          d.variable, d.value
41                          FROM UC3_Data d INNER JOIN Qualifier q
42                          ON d.file = q.file
43                          INNER JOIN Rank r
44                          ON d.variable = r.Compound AND q.expt = r.expt
45                          INNER JOIN Cmpd_info c
46                          ON d.variable = c.Compound
47                          INNER JOIN Mol_wt m
48                          ON d.variable = m.variable
49                          WHERE
50                          (q.expt = 'expt2' AND r.Confidence = 1) OR
51                          (q.expt = 'expt3' AND r.Confidence = 1)
52                          AND NOT (d.variable = 'Formate' AND d.value > 0.3)
53                          ORDER BY q.expt;")
54 dbDisconnect(con)
55
56 #Data clean up
57 clean <- data_cleanup(input, convert.NA=TRUE)
58
59 #Descriptive statistics
60 stat <- ddply(clean, 'expt', descriptive.stat, data.type='long',
61              variable=c('type', 'variable'))
62
63 #Removing compounds found only in one of Expt 2 and Expt 3
64 expt2.var <- unique(stat$variable[stat$expt=='expt2'])
65 expt3.var <- unique(stat$variable[stat$expt=='expt3'])
66
67 remove <- expt2.var[!expt2.var %in% expt3.var]
68 sub <- stat[stat$variable != remove,]
69 #General to simplest mlr model
70 check_model <- function(d, f1, f2) {
71   m1 <- lm(formula=f1, data=d)
72   m2 <- lm(formula=f2, data=d)
73   compare <- anova(m1, m2)

```

```

74   out <- compare[2,6]
75   return(out)
76 }
77
78 bonferroni <- 0.05/length(unique(sub$variable))
79 mlr1 <- ddply(sub, 'variable', mlr_custom,
80             f='value~speed+expt+process+speed:process+speed:expt')
81 mlr1$label <- ''
82 mlr1$label[mlr1[, 'Pr(>|t|)'] < bonferroni] <- '*'
83 mlr1$model <- -0.3/2
84
85 mlr2 <- ddply(sub, 'variable', mlr_custom,
86             f='value~speed+expt+process+speed:expt') #remove speed:process
87 mlr2$label <- ''
88 mlr2$label[mlr2[, 'Pr(>|t|)'] < bonferroni] <- '*'
89 mlr2$model <- -0.1/2
90
91 #compare mlr1 and mlr2
92 m1m2 <- ddply(sub, 'variable', check_model,
93             f1='value~speed+expt+process+speed:process+speed:expt',
94             f2='value~speed+expt+process+speed:expt')
95
96 mlr3 <- ddply(sub, 'variable', mlr_custom,
97             f='value~speed+expt+speed:expt') #remove process
98 mlr3$label <- ''
99 mlr3$label[mlr3[, 'Pr(>|t|)'] < bonferroni] <- '*'
100 mlr3$model <- 0.1/2
101
102 m2m3 <- ddply(sub, 'variable', check_model,
103             f1='value~speed+expt+process+speed:expt',
104             f2='value~speed+expt+speed:expt')
105
106 mlr4 <- ddply(sub, 'variable', mlr_custom,
107             f='value~speed+expt') #remove speed:expt
108 mlr4$label <- ''
109 mlr4$label[mlr4[, 'Pr(>|t|)'] < bonferroni] <- '*'
110 mlr4$model <- 0.3/2
111
112 m3m4 <- ddply(sub, 'variable', check_model,
113             f1='value~speed+expt+speed:expt',
114             f2='value~speed+expt')
115
116 #visualizing significance of terms
117 plot.data <- rbind(mlr1, mlr2, mlr3, mlr4)
118 plot.data <- plot.data[plot.data$rname != '(Intercept)',]
119 plot.data$rname <- factor(plot.data$rname, levels=c('speed', 'exptexpt3',
120                                                  'processpre-spin', 'speed:exptexpt3',
121                                                  'speed:processpre-spin'))
122 var.level <- rev(unique(plot.data$variable))
123 plot.data$variable <- factor(plot.data$variable, levels=var.level)
124
125 rname.code <- data.frame(rname=unique(plot.data$rname),
126                         rname.code=c(1,2,3,5,4))
127
128 plot.data <- merge(plot.data, rname.code, by='rname')
129 plot.data$x <- plot.data$model + plot.data$rname.code
130
131 var.code <- data.frame(variable=var.level)
132 var.code$var.code <- rev(as.numeric(var.code$variable))
133 plot.data <- merge(plot.data, var.code, by='variable')
134 plot.data$y <- plot.data$var.code
135 plot.data$var.code[plot.data$label != '*'] <- NA
136
137 plot.data$model <- factor(plot.data$model)
138
139 p <- ggplot(plot.data, aes(x=x, y=var.code)) +
140   geom_point(aes(colour=model, shape=42, size=5)) +
141   scale_x_discrete(breaks=1:5, labels=c('Speed', 'Experiment', 'Pre-Process',

```



```

142             'Speed-Experiment\nInteraction',
143             'Speed-Pre-process\nInteraction')) +
144     coord_cartesian(xlim = c(0.5, 5.5)) +
145     scale_colour_discrete(name='MLR Model', labels=c(paste('model', 1:4, sep=' '))) +
146     scale_y_continuous(breaks=var.code$var.code, labels=var.code$variable) +
147     xlab('Effect') +
148     ylab('Compound') +
149     theme_bw(11) +
150     theme(axis.ticks.x=element_blank(),
151           panel.grid.minor.x=element_blank(), panel.grid.major.x=element_blank(),
152           panel.grid.minor.y=element_blank())
153
154 ggsave('.. /Plots/Factor effects-mlr analysis/simplifying mlr-model compare2.png', p
155       , height=20, width=22, units='cm')
156
157 #Prediction from MLR
158 mlr4.predict <- ddply(sub, 'variable', mlr.predict, f='value~speed+expt')
159
160 #absolute changes
161 values <- unique(mlr4.predict[, c('speed', 'variable', 'avg', 'prediction')])
162
163 ran <- ddply(values, 'variable', mutate, r.avg1=range(avg)[1], r.avg2=range(avg)
164           [2],
165           r.pred1=range(prediction)[1], r.pred2=range(prediction)[2])
166 ran <- unique(ran)
167 ran$diff.avg <- ran$r.avg2-ran$r.avg1
168 ran$diff.pred <- ran$r.pred2-ran$r.pred1
169
170 diff <- unique(ran[ran$variable %in% c('Phenylalanine', 'Tyrosine', 'p-Cresol'),
171           c('variable', 'diff.avg', 'diff.pred')])
172
173 abs.slope <- mlr4[mlr4$variable %in% c('Phenylalanine', 'Tyrosine', 'p-Cresol'),]
174 #Percent decrease per 10 000 rpm
175 values <- unique(mlr4.predict[mlr4.predict$class=='I',
176           c('expt', 'speed', 'variable', 'avg', 'prediction')])
177
178 pred_custom <- function(d, formula){
179   m <- lm(d, formula=formula)
180   s <- summary(m)
181   out <- d
182   out$slope <- s$coefficients[2,1]
183   out$int <- s$coefficients[1,1]
184   return(out)
185 }
186
187 lm.pred <- ddply(values, c('variable', 'expt'), pred_custom,
188   formula="prediction~speed")
189
190 per.diff_custom <- function(d) {
191   d$per.diff <- NA
192   for(i in 2:nrow(d)) {
193     d$per.diff[i] <- d$slope[i]*10000/d$prediction[d$speed==10000]*100
194   }
195   return(d)
196 }
197
198 lm.pred <- ddply(lm.pred, c('variable', 'expt'), per.diff_custom)
199
200 export <- unique(lm.pred[!is.na(lm.pred$per.diff), c('expt', 'variable', 'per.diff')])
201 #write.csv(file = '../Results/Percent change per 10 000 rpm2.csv', export)
202
203 per.stat <- export[export$variable != 'p-Cresol',]
204 avg.per1 <- mean(per.stat$per.diff[per.stat$expt=='expt2'])
205 avg.per2 <- mean(per.stat$per.diff[per.stat$expt=='expt3'])
206 sd.per1 <- sd(per.stat$per.diff[per.stat$expt=='expt2'])
207 sd.per2 <- sd(per.stat$per.diff[per.stat$expt=='expt3'])
208
209 #combos of significant terms
210 speed.sig <- mlr4[mlr4$name == 'speed' & mlr4[, 'Pr(>|t|)'] < bonferroni,]
211 speed.nsig <- mlr4[mlr4$name == 'speed' & mlr4[, 'Pr(>|t|)'] >= bonferroni,]

```

```

208 expt.sig <- mlr4[mlr4$name == 'exptexpt3' & mlr4[, 'Pr(>|t|)'] < bonferroni,]
209 expt.nsig <- mlr4[mlr4$name == 'exptexpt3' & mlr4[, 'Pr(>|t|)'] >= bonferroni,]
210
211 #plotting compounds with sig bias
212 bias <- mlr4.predict[mlr4.predict$variable %in% unique(expt.sig$variable),]
213 bias$linetype <- 's'
214 bias$linetype[bias$variable %in% unique(speed.nsig$variable)] <- 'n'
215
216 p.bias <- ggplot(bias, aes(x=speed, y=value, colour=expt)) +
217   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
218   geom_point() +
219   geom_line(aes(x=speed, y=prediction, linetype=linetype)) +
220   facet_wrap(~variable, scales='free') +
221   scale_linetype_manual(values=c(5,1), guide=FALSE) +
222   scale_x_continuous(breaks=seq(0, 50000, 10000), labels=c(0:5)) +
223   scale_colour_discrete(name='Experiment', labels=c('Experiment 1', 'Experiment 2'))
224   +
225   theme_bw(11) +
226   xlab('Ultracentrifuge speed (10000 rpm)') +
227   ylab('Concentration (mM)') +
228   theme(legend.position='bottom')
229
230 ggsave('.../Plots/Factor effects-mlr analysis/mlr4-bias4.png', p.bias,
231        height=20, width=20, units='cm')
232
233 #Exploring biased compounds
234 bias.model <- mlr4[mlr4$variable %in% unique(expt.sig$variable),]
235
236 importance <- data.frame(variable=unique(bias.model$variable),
237                           intercept=bias.model$Estimate[bias.model$name=='(Intercept)'],
238                           bias=bias.model$Estimate[bias.model$name=='exptexpt3'],
239                           se=bias.model[, 'Std. Error'][bias.model$name=='exptexpt3'])
240
241 importance$smpt <- abs(importance$bias/importance$intercept*100)
242 importance$smpt.se <- abs(importance$se/importance$intercept*100)
243
244 #exporting files
245 write.csv(file='.../Results/Bias importance.csv', importance)
246 write.csv(file='.../Results/factor effect-expt sig.csv', expt.sig)
247
248 #Plotting compounds with significant slope (speed)
249 speed <- mlr4.predict[mlr4.predict$variable %in% unique(speed.sig$variable),]
250
251 p <- ggplot(speed, aes(x=speed, y=value, colour=expt)) +
252   geom_point(size=1.5) +
253   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
254   geom_line(aes(x=speed, y=prediction)) +
255   facet_wrap(~variable, scales='free', ncol=5) +
256   scale_x_continuous(breaks=seq(0, 50000, 10000), labels=c(0:5)) +
257   scale_colour_discrete(name='Experiment', labels=c('Experiment 1', 'Experiment 2'))
258   +
259   theme_bw(10) +
260   xlab('Ultracentrifuge speed (10000 rpm)') +
261   ylab('Concentration (mM)') +
262   theme(legend.position='bottom')
263
264 ggsave('.../Plots/Factor effects-mlr analysis/mlr4-speed3.png', p,
265        height=25, width=20, units='cm')
266
267 #export speed insig
268 speed.nsig <- mlr4[mlr4$name=='speed' & mlr4[, 'Pr(>|t|)'] > 0.06,]
269
270 write.csv(file='.../Results/factor effect-speed nsig.csv', speed.nsig)

```

```

2 #Ultracentrifuge Experiment 3 – Analysis of effects of processing methods
3 #ultracentrifugation vs filtration
4
5 #Pre-ambler-----
6 #Setting working directory
7 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Data source
  files/')
8
9 #Sourcing required functions
10 source('../Functions/function-data cleanup.R')
11 source('../Functions/function-descriptive statistics.R')
12
13 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 2/R
  Scripts/Functions')
14 source('../ANOVA-mean.R')
15 source('../kruskal wallis.R')
16 source('../linear model.R')
17 source('../tukey.R')
18 source('../correlation.R')
19
20 #loading required packages
21 require(RSQLite)
22 require(reshape2)
23 require(plyr)
24 require(ggplot2)
25
26 #setting working directory
27 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
  Scripts/Analysis')
28 #-----
29 #Input data
30 # Choose driver
31 drv <- dbDriver('SQLite')
32
33 # Create connection (for SQLite, this is a file)
34 con <- dbConnect(drv, '../Data source files/data.db')
35
36 expt3 <- dbGetQuery(con, "SELECT q.sampleID, q.type, q.grp, q.rep, q.class,
37                           q.speed, q.process, q.filter, q.expt, d.variable, d.value
38                           FROM UC3.Data d INNER JOIN Qualifier q
39                           ON d.file = q.file
40                           INNER JOIN Rank r
41                           ON d.variable = r.Compound AND q.expt = r.expt
42                           WHERE r.Confidence = 1 AND q.expt = 'expt3'
43                           AND NOT (d.variable = 'Formate' AND d.value > 0.3)
44                           ORDER BY q.type;")
45
46 dbDisconnect(con)
47
48 #Data clean up-----
49
50 clean <- data_cleanup(expt3, convert.NA=TRUE)
51
52 #Descriptive statistics-----
53
54 #Profile overview-----
55 plot.data <- stat
56 plot.data$label[stat$type %in% LETTERS[2:6]] <- 'ultracentrifuge'
57 plot.data$label[stat$type == 'A'] <- 'serial filter'
58 #plot.data$type <- factor(plot.data$type, levels=c('H', 'I', LETTERS[1:5]))
59
60 p <- ggplot(plot.data, aes(x=speed, y=avg, colour=label)) +
61   geom_point(aes(y=value)) +
62   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
63   facet_wrap(~variable, scale='free') +
64   scale_colour_discrete(name='Process', guide=guide_legend(reverse=TRUE)) +
65   theme_bw(14) +

```

```

66   xlab('Sample') +
67   ylab('Concentration (mM)')
68
69   ggsave('.../Plots/Processing method comparison/Profile overview.png', p,
70         height=30, width=40, units='cm')
71
72   #ANOVA-----
73
74   #filtration processes vs ultracentrifuge processes
75   anova <- ddply(stat, 'variable', aov_custom, 'class')
76   sig.anova <- anova[anova$P.val <= 0.05,]
77   write.csv(file='.../Results/anova-process.csv', sig.anova)
78
79   #tukey's HSD test on filtration vs UC processes
80   tukey.data <- stat[stat$variable %in% sig.anova$variable,]
81   tukey <- ddply(tukey.data, 'variable', tukey_custom, 'type')
82
83   bonferroni <- 0.05/length(unique(tukey.data$type))
84   sig.tukey <- tukey[tukey$P.val <= bonferroni,]
85   insig.tukey <- tukey[tukey$P.val > bonferroni,]
86
87   #See which compounds do not differ significantly between
88   #filter-processed and low speed
89   diff <- sig.tukey[sig.tukey$g1 %in% c('S','X') &
90                 sig.tukey$g2 %in% c('S','X'),]
91   same <- insig.tukey[insig.tukey$g1 %in% c('S','X') &
92                 insig.tukey$g2 %in% c('S','X'),]
93   write.csv(file='.../Results/tukey-process.csv', same)
94
95   #looking to see if filtration samples are similar to low UC speeds
96   low <- ddply(stat[stat$speed <= 30000,], 'variable', aov_custom, 'speed')
97   sig.low <- low[low$P.val <= 0.05,]
98   insig.low <- low[low$P.val > 0.05,]
99   write.csv(file='.../Results/anova-process_filt vs low speed.csv', sig.low)
100
101   #linear relationships-----
102   #seeing if linear relationships are maintained even when including
103   #filtration samples -> treating filtration samples as ultrafiltration speed 0
104   #lumping serial filtration and single-stage filtration together, since they do not
105   #differ that much
106
107   lm.result <- ddply(stat, 'variable', lm_custom, 'speed')
108   lm.sig <- lm.result[lm.result$slope.p <= 0.05,]
109   lm.result$linetype <- 'ns'
110   lm.result$linetype[lm.result$variable %in% lm.sig$variable] <- 's'
111
112   #compounds that were linear in UC samples
113   uc.lm <- read.csv('.../Results/correlation-speed.csv')
114   uc.lm <- as.character(uc.lm$variable)
115
116   #seeing if they are still linear with addition of filtration samples
117   #any uc.lm stop being linear?
118   any(uc.lm[!uc.lm %in% lm.sig$variable]) #no
119   #any compounds become linear after introduction of filter-processed samples?
120   lm.sig[!lm.sig$variable %in% uc.lm,]
121
122   corr.result <- ddply(stat, 'variable', corr_custom,
123                     x='speed', y='value', methods='pearson')
124   corr.sig <- corr.result[corr.result$P.val <= 0.05,]
125   corr.insig <- corr.result[corr.result$P.val > 0.05,]
126   write.csv(file='.../Results/correlation-process.csv', corr.sig)
127
128   #correlation of compounds that were linear in UC analysis
129   com.corr <- corr.sig[corr.sig$variable %in% uc.lm,]
130   uncom.corr <- corr.sig[!corr.sig$variable %in% uc.lm,]
131   write.csv(file='.../Results/common correlation-process.csv', com.corr)
132   write.csv(file='.../Results/uncommon correlation-process.csv', uncom.corr)
133

```

```

134 p.lm <- ggplot(plot.data, aes(x=speed, y=avg, colour=label)) +
135   geom_abline(data=lm.result, aes(slope=slope, intercept=int,
136                                   linetype=linetype), colour='grey50',) +
137   geom_point(aes(y=value), position=position_dodge(width=8000)) +
138   geom_errorbar(aes(ymax=ymax, ymin=ymin), position=position_dodge()) +
139   scale_linetype_manual(values=c(2,1)) +
140   scale_x_continuous(breaks=seq(0, 50000, 10000)) +
141   facet_wrap(~variable, scale='free', ncol=7) +
142   theme_bw(14) +
143   theme(axis.text.x=element_text(angle=-45, hjust=0, vjust=1)) +
144   xlab('Ultracentrifugation Speed (rpm)') +
145   ylab('Mean Concentration (mM)')
146
147 ggsave('~/Plots/Processing method comparison/Linear relationship.png', p.lm,
148        height=40, width=35, units='cm')

```

Centrifugal force analysis.R

```

1 #####
2 #Ultracentrifuge Experiment 3 - Analysis of effects of centrifugal force
3 #different speeds
4
5 #Pre-amble-----
6 #Setting working directory
7 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Data source
  files/')
8
9 #Sourcing required functions
10 source('~/Functions/function-data cleanup.R')
11 source('~/Functions/function-descriptive statistics.R')
12
13 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 2/R
  Scripts/Functions')
14 source('~/ANOVA-mean.R')
15 source('~/kruskal wallis.R')
16 source('~/linear model.R')
17 source('~/tukey.R')
18 source('~/correlation.R')
19
20 #loading required packages
21 require(RSQLite)
22 require(reshape2)
23 require(plyr)
24 require(ggplot2)
25
26 #setting working directory
27 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
  Scripts/Analysis')
28 #-----
29 #Input data
30 # Choose driver
31 drv <- dbDriver('SQLite')
32
33 # Create connection (for SQLite, this is a file)
34 con <- dbConnect(drv, '~/Data source files/data.db')
35
36 expt3 <- dbGetQuery(con, "SELECT q.sampleID, q.type, q.grp, q.rep, q.class,
37                            q.speed, q.process, q.filter, q.expt, c.Type, m.Weight,
38                            d.variable, d.value
39                            FROM UC3_Data d INNER JOIN Qualifier q
40                            ON d.file = q.file
41                            INNER JOIN Rank r
42                            ON d.variable = r.Compound AND q.expt = r.expt
43                            INNER JOIN Cmpd_info c
44                            ON d.variable = c.Compound
45                            INNER JOIN Mol_wt m
46                            ON d.variable = m.variable
47                            WHERE r.Confidence = 1 AND q.expt = 'expt3'")

```

```

48         AND NOT (d.variable = 'Formate' AND d.value > 0.3)
49         ORDER BY q.type;")
50
51 dbDisconnect(con)
52
53 #Data clean up-----
54 clean <- data_cleanup(expt3, convert.NA=TRUE)
55
56 #Descriptive statistics-----
57 stat <- descriptive.stat(clean, data.type='long', variable=c('type','variable'))
58
59 #profile overview-----
60 #plotting compound levels at diff UC speeds
61 p.uc <- ggplot(stat, aes(x=speed, y=avg)) +
62   geom_point(aes(y=value)) +
63   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
64   facet_wrap(~variable, scale='free') +
65   theme_bw(14) +
66   theme(axis.text.x=element_text(angle=-45, hjust=0, vjust=1)) +
67   xlab('Ultracentrifugation Speed (rpm)') +
68   ylab('Mean Concentration (mM)')
69
70 ggsave('..Plots/Speed comparison/Profile overview2.png', p.uc, height=35,width=40,
71        units='cm')
72
73 #Linear relationship-----
74 #Looking for linear relationship between concentration and speed
75
76 lm.result <- ddply(stat, 'variable', lm_custom, 'speed')
77 lm.result$linetype <- 'ns'
78 lm.result$linetype[lm.result$slope.p < 0.06] <- 's'
79
80 lm.sig <- lm.result[lm.result$slope.p < 0.06,]
81 lm.nsig <- lm.result[!lm.result$variable %in% lm.sig$variable,]
82 #write.csv(file = '../Results/lm non-significant.csv', lm.nsig)
83
84 #correlation
85 corr.result <- ddply(stat, 'variable', corr_custom,
86   x='speed', y='value', methods='pearson')
87 corr.sig <- corr.result[corr.result$p.val <= 0.05,]
88 corr.insig <- corr.result[corr.result$p.val > 0.05,]
89 #write.csv(file = '../Results/correlation-speed2.csv', corr.sig)
90 #write.csv(file = '../Results/weak correlation-speed2.csv', corr.insig)
91
92 #plotting lm-----
93 p.lm <- ggplot(stat, aes(x=speed, y=avg)) +
94   geom_point(aes(y=value)) +
95   #geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
96   stat_smooth(aes(y=value), method='lm', colour='grey85') +
97   geom_abline(data=lm.result, aes(slope=slope, intercept=int,
98     colour=linetype, linetype=linetype)) +
99   scale_linetype_manual(values=c(2,1)) +
100   scale_colour_manual(values=c('black','red')) +
101   scale_x_continuous(breaks=seq(10000, 50000, 10000)) +
102   facet_wrap(~variable, scale='free') +
103   theme_bw(14) +
104   theme(axis.text.x=element_text(angle=-45, hjust=0, vjust=1)) +
105   xlab('Ultracentrifugation Speed (rpm)') +
106   ylab('Mean Concentration (mM)')
107
108 ggsave('..Plots/Speed comparison/linear relationship4.png', p.lm,
109        height=35, width=40, units='cm')
110
111 plot.data <- stat[stat$variable %in% unique(lm.sig$variable),]
112 p.siglm <- ggplot(plot.data, aes(speed, y=avg)) +
113   geom_point(aes(y=value)) +
114   stat_smooth(aes(y=value), method='lm', colour='black') +
115   scale_x_continuous(breaks=seq(0, 50000, 10000)) +

```

```

115   facet_wrap(~variable, scale='free', ncol=6) +
116   theme_bw(14) +
117   theme(axis.text.x=element_text(angle=-45, hjust=0, vjust=1)) +
118   xlab('Ultracentrifuge Speed (rpm)') +
119   ylab('Mean Concentration (mM)')
120
121 ggsave(' ../Plots/Speed comparison/sig linear relationship2.png', p.siglm,
122        height=40, width=35, units='cm')
123
124 #Sources of variation-----
125 #Kruskal wallis test
126 kw <- ddply(stat, 'variable', kruskal_custom, 'type')
127 kw.sig <- kw[kw$p.val <= 0.05,]
128 unique(kw$variable)
129
130 #ANOVA
131 anova <- ddply(stat, 'variable', aov_custom, 'type')
132 anova.sig <- anova[anova$p.val <= 0.05,]
133 write.csv(file='../Results/uc anova.csv', anova.sig)
134
135 #Post hoc, Tukey's HSD to see where difference is coming from
136 sig.data <- stat[stat$variable %in% anova.sig$variable,]
137 tukey <- ddply(sig.data, 'variable', tukey_custom, 'type')
138 tukey.sig <- tukey[tukey$p.val < 0.05,]
139
140 plot.data <- stat[stat$variable %in% anova.sig$variable,]
141 lm.anova <- lm.result[lm.result$variable %in% anova.sig$variable,]
142
143 p.anova <- ggplot(plot.data, aes(x=speed, y=avg)) +
144   geom_point(aes(y=value)) +
145   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
146   geom_abline(data=lm.anova, aes(slope=slope, intercept=int,
147                                colour=linetype, linetype=linetype)) +
148   scale_colour_manual(values='red') +
149   scale_x_continuous(breaks=seq(0,50000,10000)) +
150   facet_wrap(~variable, scale='free', ncol=5) +
151   theme_bw(14) +
152   theme(axis.text.x=element_text(angle=-45, hjust=0, vjust=1)) +
153   xlab('Ultracentrifugation Speed (rpm)') +
154   ylab('Mean Concentration (mM)')
155
156 ggsave(' ../Plots/Speed comparison/UC-anova2.png', p.anova,
157        height=30,width=25, units='cm' )
158
159 #Range of changes as percent of avg-----
160 r <- stat
161 r <- ddply(r, 'variable', mutate, t.avg=mean(value), max=max(avg), min=min(avg))
162 r <- ddply(r, 'variable', mutate, p.max=(max-t.avg)/t.avg*100,
163          p.min=(min-t.avg)/t.avg*100)
164
165 p.range <- unique(r[,c('variable', 'Type', 'Weight', 't.avg', 'p.max', 'p.min')])
166 p.range$slope = NULL
167 p.range$int = NULL
168
169 for(i in unique(p.range$variable)) {
170   p.range$slope[p.range$variable==i] <- lm.result$slope[lm.result$variable==i]
171   p.range$int[p.range$variable==i] <- lm.result$int[lm.result$variable==i]
172 }
173
174 percent_custom <- function(d, slope='slope', int='int', x) {
175   y <- d[,slope]*x+d[,int]
176   per <- (y[2]-y[1])/y[1]*100
177
178   out <- d
179   out$perc.change <- per
180   out$abs.change <- abs(per)
181   return(out)
182 }

```

```

183
184 percent <- ddply(p.range, 'variable', percent_custom,
185                  x=seq(10000, 50000, 10000))
186
187 #write.csv(file = '../Results/Percent change per 10 000 rpm.csv', percent)
188
189 #absolute changes-----
190 values <- unique(r[,c('speed', 'variable', 'avg')])
191
192 r2 <- ddply(values, 'variable', mutate, r.max=range(avg)[2], r.min=range(avg)[1])
193 r2 <- unique(r2[,c('variable', 'r.max', 'r.min')])
194 r2$diff <- r2$r.max-r2$r.min
195
196 #control precision-----
197 ctrl <- stat[stat$speed == 0,]
198 ctrl <- ddply(ctrl, 'variable', mutate,
199               max=max(value), min=min(value))
200 ctrl <- ddply(ctrl, 'variable', mutate, p.max=(max-avg)/avg*100,
201               p.min=(min-avg)/avg*100, p.range=(max-min)/avg*100)
202 ctrl <- unique(ctrl[,c('variable', 'avg', 'se', 'max', 'min', 'p.max', 'p.min', 'p.range')
203                    ])
204 #see if percent change can be categorized into -----
205 #compound class or molecular weight
206 percent.sig <- percent[percent$variable %in% lm.sig$variable,]
207
208 sig.dec <- percent.sig[percent.sig$perc.change < 0,]
209
210 percent.stat <- data.frame(avg=mean(sig.dec$abs.change),
211                            stdev=sd(sig.dec$abs.change),
212                            max=max(sig.dec$abs.change),
213                            min=min(sig.dec$abs.change))
214
215 wt.lm <- lm(data=sig.dec, formula=Weight~abs.change)
216 wt.fit <- anova(wt.lm)
217
218 wt.corr <- cor.test(sig.dec$Weight, percent$abs.change)
219
220 p.percent <- ggplot(sig.dec, aes(x=Weight, y=abs.change)) +
221   stat_smooth(method='lm') +
222   geom_point() +
223   geom_text(aes(label=variable)) +
224   theme_bw(14) +
225   xlab('Molecular weight (g/mol)') +
226   ylab('Absolute Percent Decrease') +
227
228 ggsave('../Plots/trend by weight.png', p.percent, height=15, width=30, units='cm')
229
230 class.lm <- lm(data=percent, formula=Type~abs.change)
231
232 p.percent2 <- ggplot(sig.dec, aes(x=Type, y=abs.change)) +
233   stat_boxplot() +
234   geom_point() +
235   theme_bw(14) +
236   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))

```

Cross-experiment analysis.R

```

1 #####
2 //
3 #Comparison of profiles from Expt 3, Expt 1, Expt 2 and preliminary expt
4
5 #pre-amble-----
6 #Setting working directory
7 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
8
9 #Sourcing required functions
10 source('function-data cleanup.R')

```



```

10 source('function-descriptive_statistics.R')
11
12 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 2/R
    Scripts/Functions')
13 source('./ANOVA-mean.R')
14 source('./kruskal wallis.R')
15 source('./linear_model.R')
16 source('./tukey.R')
17 source('./correlation.R')
18
19 #loading required packages
20 require(RSQLite)
21 require(reshape2)
22 require(plyr)
23 require(ggplot2)
24
25 #setting working directory
26 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
    Scripts/Analysis')
27 source('../Functions/multiple linear regression.R')
28 source('../Functions/mlr_prediction.R')
29 #
30 #Input data
31 # Choose driver
32 drv <- dbDriver('SQLite')
33
34 # Create connection (for SQLite, this is a file)
35 con <- dbConnect(drv, '../Data source files/data.db')
36
37 input <- dbGetQuery(con, "SELECT q.SampleID, q.type, q.grp, q.rep, q.class,
38                             q.speed, q.process, q.filter, q.expt, d.variable, d.value
39                             FROM UC3_Data d INNER JOIN Qualifier q
40                             ON d.file = q.file
41                             INNER JOIN Rank r
42                             ON d.variable = r.Compound AND q.expt = r.expt
43                             WHERE
44                             (q.expt == 'expt1' AND q.class != 'II' AND r.Confidence = 1) OR
45                             (q.expt == 'expt2' AND q.class != 'II' AND r.Confidence = 1) OR
46                             q.expt IN ('expt3', 'preliminary') AND r.Confidence = 1
47                             AND NOT (d.variable = 'Formate' AND d.value > 0.3)
48                             ORDER BY q.type;")
49
50 input2 <- dbGetQuery(con, "SELECT q.SampleID, q.type, q.grp, q.rep, q.class,
51                             q.speed, q.process, q.filter, q.expt, d.variable, d.value
52                             FROM UC3_Data d INNER JOIN Qualifier q
53                             ON d.file = q.file
54                             INNER JOIN Rank r
55                             ON d.variable = r.Compound AND q.expt = r.expt
56                             WHERE
57                             (q.expt == 'expt2' AND q.class != 'II' AND r.Confidence = 1) OR
58                             q.expt = 'expt3' AND r.Confidence = 1
59                             AND NOT (d.variable = 'Formate' AND d.value > 0.3)
60                             ORDER BY q.type;")
61
62 dbDisconnect(con)
63
64
65 #Data clean up
66 clean <- data_cleanup(input2, convert.NA=TRUE)
67
68 #Descriptive statistics
69 stat <- ddply(clean, 'expt', descriptive.stat, data.type='long',
70               variable=c('type', 'variable'))
71
72 #profile overview
73 p.uc2 <- ggplot(stat, aes(x=speed, y=avg, colour=expt)) +
74   geom_point(aes(y=value), alpha=0.8) +
75   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +

```

```

76 facet_wrap(~variable, scale='free') +
77 theme_bw(14) +
78 theme(axis.text.x=element_text(angle=-45, hjust=0, vjust=1),
79 legend.position='bottom') +
80 xlab('Ultracentrifugation Speed (rpm)') +
81 ylab('Mean Concentration (mM)')
82
83 ggsave('.. /Plots/Cross-experiment comparison/Combined expt-expt2, expt3-2.png',
84 p.uc2, height=35,width=35, units='cm')
85
86 #Removing compounds found only in one of Expt 2 and Expt 3-----
87 expt2.var <- unique(stat$variable[stat$expt=='expt2'])
88 expt3.var <- unique(stat$variable[stat$expt=='expt3'])
89
90 remove <- expt2.var[!expt2.var %in% expt3.var]
91 sub <- stat[stat$variable != remove,]
92
93 #Getting equation for multiple linear regression plane-----
94 #checking for bias and interaction
95 mlr <- ddply(sub, 'variable', mlr_custom, f='value~speed+expt+speed:expt')
96
97 #compounds with insignificant slope
98 nsig.slope <- mlr[mlr[, 'Pr(>|t|)'] > 0.05 & mlr$name == 'speed',]
99
100 #removing compounds with insignificant slope (don't need to be plotted)
101 mlr.sub <- mlr[!mlr$variable %in% unique(nsig.slope$variable),]
102
103 #getting compounds with only significant bias, interaction or both
104 s.bias <- unique(mlr.sub$variable[mlr.sub[, 'Pr(>|t|)'] < 0.06
105 & mlr.sub$name == 'exptexpt3'])
106 n.intxn <- unique(mlr.sub$variable[mlr.sub[, 'Pr(>|t|)'] > 0.06
107 & mlr.sub$name == 'speed:exptexpt3'])
108 n.bias <- unique(mlr.sub$variable[mlr.sub[, 'Pr(>|t|)'] > 0.06
109 & mlr.sub$name == 'exptexpt3'])
110 s.intxn <- unique(mlr.sub$variable[mlr.sub[, 'Pr(>|t|)'] < 0.06
111 & mlr.sub$name == 'speed:exptexpt3'])
112
113 bias <- unique(mlr.sub$variable[mlr.sub$variable %in% s.bias
114 & mlr.sub$variable %in% n.intxn])
115 intxn <- unique(mlr.sub$variable[mlr.sub$variable %in% n.bias
116 & mlr.sub$variable %in% s.intxn])
117 both <- unique(mlr.sub$variable[mlr.sub$variable %in% s.bias
118 & mlr.sub$variable %in% s.intxn])
119 neither <- unique(mlr.sub$variable[mlr.sub$variable %in% n.bias
120 & mlr.sub$variable %in% n.intxn])
121
122 #Exporting compounds that followed same trend in both experiments
123 export.neither <- mlr.sub[mlr.sub$variable %in% neither, c(1,5,6)]
124 export.neither <- dcast(export.neither, variable~name, value.var='Pr(>|t|)')
125 write.csv(file='.. /Results/mlr_neither.csv', export.neither)
126
127 #predicting values along the multiple linear regression plane-----
128 mlr.predict <- ddply(sub, 'variable', mlr_predict, f='value~speed+expt+speed:expt')
129
130 #visualizing bias and interaction-----
131 plot.data <- mlr.predict
132 plot.data.bias <- plot.data[plot.data$variable %in% bias,]
133 plot.data.intxn <- plot.data[plot.data$variable %in% intxn,]
134 plot.data.both <- plot.data[plot.data$variable %in% both,]
135 plot.data.neither <- plot.data[plot.data$variable %in% neither,]
136
137 p <- ggplot(plot.data.both, aes(x=speed, y=value, colour=expt)) +
138 geom_point() +
139 geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
140 geom_line(aes(x=speed, y=prediction)) +
141 facet_wrap(~variable, scales='free') +
142 scale_x_continuous(breaks=seq(0, 50000, 10000)) +
143 scale_colour_discrete(name='Experiment', labels=c('Experiment 1', 'Experiment 2'))

```

```

+
144 theme_bw(11) +
145 xlab('Ultracentrifuge speed (rpm)') +
146 ylab('Concentration (mM)') +
147 theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
148
149 ggsave('~/Plots/Cross-experiment comparison/multiple linear regression/mlr bias
and interaction2.png', p,
150 height=10.5, width=18, units='cm')

```

Cross-experiment analysis-pca.R

```

1 #####
2 #Cross-experiment analysis-pca
3 #Comparison of profiles from Expt 1, Expt 2 and preliminary expt
4
5 #pre-ambles-----
6
7 #loading required packages
8 require(RSQLite)
9 require(ggplot2)
10
11 #sourcing required functions
12 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions/')
13 source('function-pca.R')
14 source('function-plotting 2D score.R')
15 source('function-plotting 2D loading.R')
16 source('function-pca diagnostics.R')
17 source('function-data cleanup.R')
18 source('function-scaling.R')
19
20
21 #setting working directory
22 setwd('C:/Users/Sandi/Dropbox/Projects/Ultracentrifugation Samples/Experiment 3/R
Scripts/Analysis')
23 #
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

```

```

24 #Input data
25 # Choose driver
26 drv <- dbDriver('SQLite')
27
28 # Create connection (for SQLite, this is a file)
29 con <- dbConnect(drv, '../Data source files/data.db')
30
31 input <- dbGetQuery(con, "SELECT q.sampleID, q.type, q.grp, q.rep, q.class,
32 q.speed, q.process, q.filter, q.expt, d.variable, d.value
33 FROM UC3_Data d INNER JOIN Qualifier q
34 ON d.file = q.file
35 INNER JOIN Rank r
36 ON d.variable = r.Compound AND q.expt = r.expt
37 WHERE
38 (q.expt = 'expt1' AND q.class != 'II' AND r.Confidence = 1) OR
39 (q.expt = 'expt2' AND q.class != 'II' AND r.Confidence = 1) OR
40 q.expt IN ('expt3', 'preliminary') AND r.Confidence = 1
41 AND NOT (d.variable = 'Formate' AND d.value > 0.3)
42 ORDER BY q.type;")
43
44 input2 <- dbGetQuery(con, "SELECT q.SampleID, q.type, q.grp, q.rep, q.class,
45 q.speed, q.process, q.filter, q.expt, d.variable, d.value
46 FROM UC3_Data d INNER JOIN Qualifier q
47 ON d.file = q.file
48 INNER JOIN Rank r
49 ON d.variable = r.Compound AND q.expt = r.expt
50 WHERE
51 (q.expt = 'expt2' AND q.class != 'II' AND r.Confidence = 1) OR
52 q.expt = 'expt3' AND r.Confidence = 1
53 AND NOT (d.variable = 'Formate' AND d.value > 0.3)

```

```

54         ORDER BY q.type;")
55
56 expt3 <- dbGetQuery(con, "SELECT q.SampleID, q.type, q.grp, q.rep, q.class,
57       q.speed, q.process, q.filter, q.expt, d.variable, d.value
58       FROM UC3_Data d INNER JOIN Qualifier q
59       ON d.file = q.file
60       INNER JOIN Rank r
61       ON d.variable = r.Compound AND q.expt = r.expt
62       WHERE
63       q.expt = 'expt3' AND r.Confidence = 1
64       AND NOT (d.variable = 'Formate' AND d.value > 0.3)
65       ORDER BY q.type;")
66 dbDisconnect(con)
67
68
69 #Data clean up-----
70 clean <- data_cleanup(expt3, convert.NA=FALSE)
71
72 #performing the PCA-----
73 pca.data = perform_pca(clean, sampleID='sampleID',
74       variable='variable', value='value',
75       form.y=c('sampleID','type','grp','rep','class',
76       'speed','process','filter','expt'),
77       form.x='variable',
78       scale='UV', center=TRUE, convert.NA=TRUE)
79
80 #performing diagnostics on PCA-----
81 diagnostic.results = scree_plot(pca.data)
82
83 #unpacking the diagnostic results
84 diagnostic.number = diagnostic.results[['number']]
85 diagnostic.explained = diagnostic.results[['explained']]
86 diagnostic.plot = diagnostic.results[['plot']]
87 diagnostic.plot = diagnostic.plot + ggtitle('Scree Plot (UVN)')
88
89 #printing out diagnostic results
90 print(diagnostic.number)
91 print(diagnostic.explained)
92 print(diagnostic.plot)
93
94 #plotting PCA-processed data-----
95 #Score plot 2D
96 # Note: Arguments for score plot function are
97 # (pca_data, qualifier_data, PCs=1:4, colour=NULL, shape=NULL, size=NULL,
98 # label=NULL, title=NULL, legend.ori='horizontal', black-white=FALSE)
99
100 score_plot = score_plot_2d(pca.data, PCs=1:4,
101       colour='speed', shape='expt',
102       legend.ori='vertical',
103       black-white=FALSE)
104
105 p.score <- score_plot[['p']]
106 collected.data <- score_plot[['collected.data']]
107
108 p.score = p.score +
109   #scale_shape_discrete(name='Process condition',
110   # labels=c('Pre-spun', 'UC-processed',
111   # 'Not pre-spun', 'UC-processed', 'Filter-processed'))
112   +
113   scale_shape_discrete(name='Pre-process treatment', labels=c('Pre-spun', 'Not pre-
114   spun')) +
115   #scale_shape_discrete(name='Experiment', labels=c('Experiment 1', 'Experiment 2', '
116   Preliminary')) +
117   guides(colour=guide_colourbar('Speed (rpm)')) +
118   theme_bw(14)
119
120 ggsave('.. /Plots/Cross-experiment comparison/pca/score-expt3.png', p.score,
121       height=10, width=25, units='cm')

```

C.3 Effects of Varying Inoculum Biomass Proportions

Data source.R

```

1 #####
2 #Creating database for Biomass Samples
3
4 #sourcing required functions
5 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
6 source('function-summary preprocessing.R')
7
8 #loading required packages
9 require(reshape2)
10 require(RSQLite)
11
12 #Setting working directory
13 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass samples/R Scripts/Data source files'
14 )
15 #preprocessing raw data
16 raw <- summary.preprocess('concentrations2.csv')
17
18 data <- raw[['data']]
19
20 melt.data <- melt(data)
21 melt.data$value <- melt.data$value/0.9
22
23 #Writing qualifiers
24 #reading sample names and scan order
25 sample.list <- read.csv('Sample list.csv')
26 scan <- read.csv('Scan order.csv')
27
28 #defining profiles according to scan order
29 qualifier <- data.frame(file=data$file, sampleID=scan$sampleID)
30
31 #re-ordering by sample ID
32 qualifier <- qualifier[order(qualifier$sampleID),]
33
34 #adding sample_name to qualifier
35 qualifier$sample_name <- sample.list$sample_name
36
37 qualifier$run <- qualifier$sample_name
38 qualifier$run <- gsub('.*(?=Run )([0-9]){2}.*', '\\1', qualifier$run, perl=TRUE)
39
40 qualifier$vessel <- qualifier$sample_name
41 qualifier$vessel <- gsub('.*(?=V )([0-9]){2}.*', '\\1', qualifier$vessel, perl=TRUE)
42
43 qualifier$donor <- qualifier$sample_name
44 qualifier$donor <- gsub('.*(?=, )([A-Z0-9]){2,4}.*', '\\1',
45                       qualifier$donor, perl=TRUE)
46
47 qualifier$day <- qualifier$sample_name
48 qualifier$day <- gsub('.*(?=Day )([0-9]){1,2}.*', '\\1', qualifier$day, perl=TRUE)
49 qualifier$day[qualifier$day == 'Media control'] <- 0
50
51 #Compound Ranking
52 rank <- read.csv('Compound ranking.csv', encoding='UTF-8')
53 rank <- rank[,1:7]
54
55 #Compound info
56 #metabolite pathway
57 path.data <- read.csv('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/
58                       Data source files/Compound pathway annotation.csv',
59                       header=FALSE)
60 path.data <- path.data[1:153,]
61
62 #Getting compound types
63 cmpd.type <- data.frame(Compound=t(path.data[1,2:ncol(path.data)]),

```

```

63                                     Type=t(path.data[2,2:ncol(path.data)]))
64 colnames(cmpd.type) <- c('Compound', 'Type')
65
66 # Create the connection (no file needed) -----
67
68 # Choose driver
69 drv = dbDriver('SQLite')
70
71 # Create connection (for SQLite, this is a file)
72 con = dbConnect(drv, 'data.db')
73
74 # Creating database tables
75 data.table <- dbWriteTable(con, 'Data', melt.data, overwrite=TRUE)
76
77 qualifier.table <- dbWriteTable(con, 'Qualifier', qualifier, overwrite=TRUE)
78
79 rank.table <- dbWriteTable(con, 'Rank', rank, overwrite=TRUE)
80
81 cmpd.table <- dbWriteTable(con, 'Cmpd-info', cmpd.type, overwrite=TRUE)
82
83 tables <- dbGetQuery(con, "SELECT * FROM sqlite_master WHERE type='table';")
84 print(tables)
85
86 dbDisconnect(con)

```

profiling variability analysis.R

```

1 #####
2 #####
3 #Profiling variability anlaysis
4
5 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
6 source('function-descriptive statistics.R')
7 source('function-data cleanup.R')
8
9 #loading required packages
10 require(RSQLite)
11 require(plyr)
12
13 #setting working directory
14 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass Samples/R Scripts/Analysis')
15
16 #Input data-----
17 #Choose driver
18 drv <- dbDriver('SQLite')
19
20 #Create connecction (for SQLite, this is a file)
21 con <- dbConnect(drv, '../Data source files/data.db')
22
23 input <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
24                             d.variable, d.value
25                             FROM Data d INNER JOIN Qualifier q
26                             ON d.file=q.file
27                             INNER JOIN Rank r
28                             ON d.variable = r.Compound
29                             WHERE r.Confidence = 1
30                             AND q.sampleID IN ('H1', 'H2', 'H3', 'R1', 'R2', 'R3')
31                             ORDER BY q.sampleID;")
32
33 dbDisconnect(con)
34
35 #Data clean up-----
36 #Remove DSS, contminants and compounds of zero concentration
37 clean.data <- data_cleanup(input, 'sampleID', 'variable', 'value',
38                             form.y=c('sampleID', 'run', 'vessel', 'donor', 'day'),
39                             form.x='variable', contaminants=NULL, export=FALSE,
40                             convert.NA=TRUE)

```

```

41 #Descriptive statistics-----
42 #first assigning general sampleID
43 clean.data$set <- clean.data$sampleID
44 clean.data$set <- gsub('^[A-Z])(?=[0-9]).*', '\\1', clean.data$set, perl=TRUE)
45
46 stat <- ddply(clean.data, 'set', descriptive.stat, data.type='long',
47           variable='variable',)
48
49 #Residuals-----
50 stat$res <- stat$value-stat$avg
51
52 #Studentized residual
53 stat$sres <- abs(stat$res)/stat$stdev
54
55 #plot(density(stat$sres))
56
57 #getting sum of residuals for variable
58 custom_residual <- function(data) {
59   out <- ddply(data, 'variable', mutate, sum.res=sum(res))
60   return(out)
61 }
62
63 stat <- ddply(stat, 'set', custom_residual)
64
65 #plot(density(stat$sum.res))
66
67 #Range-----
68 #getting range of variable concentrations
69 custom_range <- function(data){
70   d <- data
71   out <- ddply(d, 'variable', mutate, min=min(value), max=max(value))
72   out$range <- out$min-out$max
73   return(out)
74
75 }
76 stat <- ddply(stat, 'set', custom_range)
77
78 #range as percent of mean
79 stat$prange <- abs(stat$range)/stat$avg*100
80 plot(density(stat$prange))
81
82 png(filename='../Results/Percent_range2.png',
83      width=10, height=7, units='cm', res=240)
84 plot(density(stat$prange), main='')
85 dev.off()
86
87 prange <- unique(stat[, c(6,8,ncol(stat))])
88 lrange <- prange[prange$prange >= 25,]
89 lrange <- lrange[order(lrange$prange),]
90 write.csv(file='../Results/Percent_range2.csv', x=lrange)
91
92 #Relative standard deviation-----
93 stat$rsd <- stat$stdev/stat$avg*100
94
95 plot(density(stat$rsd, bw=1))
96
97 png(filename='../Plots/precision/RSD.png',
98      width=10, height=7, units='cm', res=240)
99 plot(density(stat$rsd, bw=1), main='')
100 dev.off()
101
102 high.rsd <- unique(stat[stat$rsd > 15,c('set','variable','rsd')])
103 write.csv(file='../Results/rsd.csv',high.rsd)
104 #Compounds with values greater than 2*sd-----
105 stat$s2 <- 2*stat$stdev
106
107 sd1 <- stat[abs(stat$res) > stat$stdev,]
108 sd2 <- stat[abs(stat$res) > stat$s2,]

```

```

109
110 sub.stat <- unique(stat[,c(8,2:6,9:ncol(stat))])
111
112 #Precision anlaysis ANOVA
113 aov_custom <- function(data, group) {
114   d <- data
115
116   grand <- mean(d$value) #grand mean
117
118   ss.grp <- ddply(d, group, mutate, avg=mean(value)) #within-group mean
119   ss.grp$grand <- grand #adding grand mean
120
121   ssr <- sum((ss.grp$value - ss.grp$avg)**2) #sum of sqr of within groups
122   ssa <- sum((ss.grp$avg - ss.grp$grand)**2) #sum of sqr of between groups
123
124   dfr <- nrow(d)-length(unique(d[,group])) #n-q
125   dfa <- length(unique(d[,group]))-1 #q-1
126
127   msa <- ssa/dfa #mean square between group
128   msr <- ssr/dfr #mean square within group
129
130   f <- msa/msr #f statistic
131   p.val <- 1-pf(f, dfa, dfr)
132   out <- data.frame(variable=unique(d$variable), f.val=f, p.val=p.val)
133   return(out)
134 }
135
136 anova <- ddply(clean.data, 'variable', aov_custom, 'set')
137
138 s.anova <- anova[anova$p.val <=0.05,]
139
140 write.csv(file='../Results/ANOVA-triplicates2.csv', s.anova)
141 #plot(density(anova$f.val))
142 p <- ggplot(stat, aes(x=set, y=value)) +
143   geom_point() +
144   stat_boxplot(aes(ymax=ymax, ymin=ymin)) +
145   facet_wrap(~variable, scales='free') +
146   theme_bw(14)
147
148 #Student's t-test
149 custom_t <- function(data, group){
150   grps <- unique(data[,group])
151   g1 <- data[data[,group] == grps[1],]
152   g2 <- data[data[,group] == grps[2],]
153
154   t <- t.test(g1$value, g2$value)
155   out <- data.frame(statistic=t$statistic, p.val=t$p.value)
156   return(out)
157 }
158
159 t <- ddply(stat, 'variable', custom_t, 'set')
160 s.t <- t[t$p.val <= 0.05,]
161
162 write.csv(file='../Results/precision-t.csv', s.t)

```

duplicate variability analysis.R

```

1 #####
2 #Duplicate run variability analysis
3
4 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
5 source('function-descriptive statistics.R')
6 source('function-data cleanup.R')
7
8 #loading required packages
9 require(RSQLite)
10 require(plyr)
11 require(ggplot2)

```



```

12
13 #setting working directory
14 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass Samples/R Scripts/Analysis')
15
16 #Input data-----
17 #Choose driver
18 drv <- dbDriver('SQLite')
19
20 #Create connection (for SQLite, this is a file)
21 con <- dbConnect(drv, '../Data source files/data.db')
22
23 input <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
24                           d.variable, d.value
25                           FROM Data d INNER JOIN Qualifier q
26                           ON d.file=q.file
27                           INNER JOIN Rank r
28                           ON d.variable = r.Compound
29                           WHERE r.Confidence = 1 AND q.run != 'Media control'
30                           ORDER BY q.sampleID;")
31
32 input1 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
33                             d.variable, d.value
34                             FROM Data d INNER JOIN Qualifier q
35                             ON d.file=q.file
36                             INNER JOIN Rank r
37                             ON d.variable = r.Compound
38                             WHERE r.Confidence = 1
39                             AND q.run IN ('30', '33')
40                             ORDER BY q.sampleID;")
41
42 input2 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
43                             d.variable, d.value
44                             FROM Data d INNER JOIN Qualifier q
45                             ON d.file=q.file
46                             INNER JOIN Rank r
47                             ON d.variable = r.Compound
48                             WHERE r.Confidence = 1
49                             AND q.run IN ('31', '32')
50                             ORDER BY q.sampleID;")
51
52 dbDisconnect(con)
53
54 #Data clean up-----
55 #Remove DSS, contaminants and compounds of zero concentration
56 clean.data <- data_cleanup(input, 'sampleID', 'variable', 'value',
57                             form.y=c('sampleID', 'run', 'vessel', 'donor', 'day'),
58                             form.x='variable', contaminants=NULL, export=FALSE,
59                             convert.NA=TRUE)
60
61 #Calculating mean values of triplicates in set2 and input-----
62 if(any(clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3'))){
63   H.rep <- clean.data[clean.data$sampleID %in% c('H1', 'H2', 'H3'),]
64   R.rep <- clean.data[clean.data$sampleID %in% c('R1', 'R2', 'R3'),]
65
66   H <- ddply(H.rep, 'variable', mutate, avg=mean(value))
67   H$sampleID <- 'H'
68   H <- unique(H[, -7])
69   colnames(H)[7] <- 'value'
70
71   R <- ddply(R.rep, 'variable', mutate, avg=mean(value))
72   R$sampleID <- 'R'
73   R <- unique(R[, -7])
74   colnames(R)[7] <- 'value'
75
76   #removing H1, H2, H3, R1, R2, R3
77   clean.data <- clean.data[!clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3'),]
78   #adding mean H and R

```

```

79 clean.data <- rbind(clean.data, H, R)
80 print('here')
81 }
82
83
84 #casting data to give non-profiled compounds concentration of zero
85 data.cast <- dcast(clean.data, sampleID+run+vessel+donor+day~variable,
86 value.var='value')
87
88 data.cast[is.na(data.cast)] <- 0
89
90 #melting back to work with
91 data.melt <- melt(data.cast, id.vars=c('sampleID','run','vessel','donor','day'))
92
93 #adding ID for donor-run combo
94 data.melt$donor.run <- paste(data.melt$donor, data.melt$run, sep='.')
95
96 #Descriptive statistics
97 #stat <- descriptive.stat(data.melt, data.type='long', variable='variable')
98
99 stat <- ddply(data.melt, c('variable','donor','run'), descriptive.stat,
100 data.type='long', variable='variable')
101
102 stat <- ddply(stat, c('variable','donor'), mutate, run.avg=mean(value))
103 stat$per.diff <- abs((stat$avg-stat$run.avg))/stat$run.avg*100
104
105 #ANOVA
106 aov_custom <- function(data, group) {
107   d <- data
108
109   grand <- mean(d$value) #grand mean
110
111   ss.grp <- ddply(d, group, mutate, avg=mean(value)) #within-group mean
112   ss.grp$grand <- grand #adding grand mean
113
114   ssr <- sum((ss.grp$value - ss.grp$avg)**2) #sum of sqr of within groups
115   ssa <- sum((ss.grp$avg - ss.grp$grand)**2) #sum of sqr of between groups
116
117   dfr <- nrow(d)-length(unique(d[,group])) #n-q
118   dfa <- length(unique(d[,group]))-1 #q-1
119
120   msa <- ssa/dfa #mean square between group
121   msr <- ssr/dfr #mean square within group
122
123   f <- msa/msr #f statistic
124   p.val <- 1-pf(f, dfa, dfr)
125   #out <- data.frame(variable=unique(d$variable), f.val=f, p.val=p.val)
126   out <- d
127   out$f.val <- f
128   out$p.val <- p.val
129   return(out)
130 }
131
132 anova <- ddply(stat, c('variable','donor'), aov_custom, group='run')
133 s.anova <- anova[anova$p.val <= 0.05 & !is.na(anova$p.val),]
134 #write.csv(file='../Results/Run comparison-31 vs 32.csv', s.anova)
135
136 #Tukey with bonferroni correction
137 tukey_custom <- function(data, group) {
138   grps <- unique(data[,group])
139   pair <- combn(grps,2)
140   out <- c()
141   for(i in 1:ncol(pair)) {
142     grp1 <- data[data[,group] == pair[1,i],]
143     grp2 <- data[data[,group] == pair[2,i],]
144     grp <- rbind(grp1, grp2) #both groups being compared
145
146     grand <- mean(data$value) #grand mean

```

```

147 ss.grp <- ddply(grp, group, mutate, avg=mean(value)) #within-group mean
148 ss.grp$grand <- grand #adding grand mean
149 ssr <- sum((ss.grp$value - ss.grp$avg)**2) #sum of sqr of within groups
150 dfr <- nrow(data)-length(unique(data[,group])) #n-q; n=# of observations,
151 #q= # of groups
152 msr <- ssr/dfr
153
154 avg1 <- mean(grp1$value)
155 avg2 <- mean(grp2$value)
156
157 n <- nrow(grp1) #number of observations in a group (this is parametric test)
158 q <- length(unique(data[,group])) #number of groups
159 se <- sqrt(msr/n)
160 q.stat <- (max(c(avg1, avg2))-min(c(avg1, avg2)))/se
161 p.val <- ptukey(q.stat, q, df=dfr, lower.tail=FALSE)
162
163 #print(pair[,i])
164 #print(grp1$value)
165 #print(grp2$value)
166 #print(msr)
167 #print(se)
168 #print(p.val)
169 #cat('\n\n')
170 output <- data.frame(g1=pair[1,i], g2=pair[2,i], p.val=p.val)
171
172 out <- rbind(out, output)
173
174 }
175 return(out)
176 }
177
178 tukey <- ddply(data.melt, 'variable', tukey_custom, group='donor.run')
179 bonferroni <- 0.05/6
180 s.tukey <- tukey[tukey$p.val <= bonferroni & !is.na(tukey$p.val),]
181
182 #t-test
183 t_custom <- function(data, group) {
184
185   grps <- unique(data[,group])
186
187   grp1 <- data[data[,group] == grps[1],]
188   grp2 <- data[data[,group] == grps[2],]
189   grp <- rbind(grp1, grp2) #both groups being compared
190
191   result <- t.test(grp1$value, grp2$value)
192   p.val <- result$p.value
193   statistic <- result$statistic
194   output <- data.frame(g1=unique(grp1[,group]), g2=unique(grp2[,group]),
195                       stat=statistic, p.val=p.val)
196
197   return(output)
198 }
199
200 t.result <- ddply(stat,c('variable','donor'), t_custom, group='run')
201 t.sig <- t.result[t.result$p.val <= 0.05 & !is.na(t.result$p.val),]
202
203 bonferroni <- 0.05/sqrt(length(unique(stat$variable))+length(unique(stat$donor)))
204 t.sig2 <- t.result[t.result$p.val <= bonferroni & !is.na(t.result$p.val),]
205
206 for(i in c('RP','RRP','RP2','RRP2')) {
207   d1 <- unique(stat[stat$donor==i,c('variable','donor','per.diff')])
208   d2 <- t.sig[t.sig$donor == i,]
209   export <- merge(d2,d1,by='variable')
210   export <- unique(export[,c('donor.y','variable','stat','p.val','per.diff')])
211   write.csv(file=paste('../Results/run variability-',i,'.csv',sep=''),export)
212 }
213 write.csv(file='../Results/run variability-t_bonferroni.csv', t.sig2)
214

```

```

215 stat.sig <- stat[stat$variable %in% unique(t.sig$variable),]
216 stat.sig2 <- stat[stat$variable %in% unique(t.sig2$variable),]
217
218 t.list <- list(rp=unique(t.sig$variable[t.sig$donor=='RP']),
219               rrp=unique(t.sig$variable[t.sig$donor=='RRP']),
220               rp2=unique(t.sig$variable[t.sig$donor=='RP2']),
221               rrp2=unique(t.sig$variable[t.sig$donor=='RRP2']))
222
223
224 #Boxplot
225 #significant by t test
226 donor <- c('RP', 'RRP', 'RP2', 'RRP2')
227 height <- c(15, 10, 8, 10)
228 width <- c(20, 20, 15, 20)
229 n.col <- c(4, 4, 3, 4)
230 for(i in 1:length(donor)){
231   pbox.data <- stat[stat$donor == donor[i] & stat$variable %in% t.list[[i]],]
232
233   p.box <- ggplot(pbox.data, aes(x=run, y=value)) +
234     stat_boxplot() +
235     geom_point() +
236     facet_wrap(~variable, scales='free', ncol=n.col[i]) +
237     theme_bw(14) +
238     xlab('Run') +
239     ylab('Concentration (mM)')
240
241   #ggsave('.../Plots/Run variability/anova2-sig-30 vs 33.png', p.box, height=15,
242           width=25, units='cm')
243
244   ggsave(paste('.../Plots/Run variability/anova', i, '.png'), p.box,
245           height=height[i], width=width[i], units='cm')
246 }

```

pca analysis.R

```

1 #####
2 #PCA Analysis of Biomass samples
3
4 #sourcing required functions
5 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
6 source('function-data cleanup.R')
7 source('function-pca.R')
8 source('function-pca diagnostics.R')
9 source('function-plotting 2D score.R')
10 source('function-plotting 2D loading.R')
11
12 #loading required packages
13 require(RSQLite)
14 require(plyr)
15
16 #setting working directory
17 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass Samples/R Scripts/Analysis')
18
19 #Input data
20 #Choose driver
21 drv <- dbDriver('SQLite')
22
23 #Create connection (for SQLite, this is a file)
24 con <- dbConnect(drv, '../Data source files/data.db')
25
26 input <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
27                           d.variable, d.value
28                           FROM Data d INNER JOIN Qualifier q
29                           ON d.file=q.file
30                           INNER JOIN Rank r
31                           ON d.variable = r.Compound
32                           WHERE r.Confidence = 1
33                           ORDER BY q.sampleID;")

```

```

34
35 input2 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
36                             d.variable, d.value
37                             FROM Data d INNER JOIN Qualifier q
38                             ON d.file=q.file
39                             INNER JOIN Rank r
40                             ON d.variable = r.Compound
41                             WHERE r.Confidence = 1 AND q.run != 'Media control'
42                             ORDER BY q.sampleID;")
43
44 set1 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
45                             d.variable, d.value
46                             FROM Data d INNER JOIN Qualifier q
47                             ON d.file=q.file
48                             INNER JOIN Rank r
49                             ON d.variable = r.Compound
50                             WHERE r.Confidence = 1 AND q.run IN (30,33)
51                             ORDER BY q.sampleID;")
52
53 set2 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
54                             d.variable, d.value
55                             FROM Data d INNER JOIN Qualifier q
56                             ON d.file=q.file
57                             INNER JOIN Rank r
58                             ON d.variable = r.Compound
59                             WHERE r.Confidence = 1 AND q.run IN (31,32)
60                             ORDER BY q.sampleID;")
61
62 set3 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
63                             d.variable, d.value
64                             FROM Data d INNER JOIN Qualifier q
65                             ON d.file=q.file
66                             INNER JOIN Rank r
67                             ON d.variable = r.Compound
68                             WHERE r.Confidence = 1 AND q.sampleID IN ('H1','H2','H3',
69                             'R1','R2','R3')
70                             ORDER BY q.sampleID;")
71
72 dbDisconnect(con)
73
74 #Data clean up
75 #Remove DSS, contaminants and compounds of zero concentration
76 clean.data <- data.cleanup(set3, 'sampleID', 'variable', 'value',
77                             form.y=c('sampleID','run','vessel','donor','day'),
78                             form.x='variable', contaminants=NULL, export=FALSE,
79                             convert.NA=FALSE)
80 if(FALSE){
81   export.data <- dcast(clean.data, formula=sampleID+run+vessel+donor+day~variable,
82                         value.var='value')
83   write.csv(file='../metabolite_data2.csv',export.data)
84 }
85
86 #Calculating mean values of triplicates in set2 and input
87 if(any(clean.data$sampleID %in% c('H1','H2','H3','R1','R2','R3'))){
88   H.rep <- clean.data[clean.data$sampleID %in% c('H1','H2','H3'),]
89   R.rep <- clean.data[clean.data$sampleID %in% c('R1','R2','R3'),]
90
91   H <- ddply(H.rep, 'variable', mutate, avg=mean(value))
92   H$sampleID <- 'H'
93   H <- unique(H[, -7])
94   colnames(H)[7] <- 'value'
95
96   R <- ddply(R.rep, 'variable', mutate, avg=mean(value))
97   R$sampleID <- 'R'
98   R <- unique(R[, -7])
99   colnames(R)[7] <- 'value'
100
101 #removing H1,H2,H3,R1,R2,R3

```

```

101   clean.data <- clean.data[!clean.data$sampleID %in% c('H1','H2','H3','R1','R2','R3
102   '),]
103   #adding mean H and R
104   clean.data <- rbind(clean.data, H, R)
105   print('here')
106 }
107 #adding columns of coded run and donor labels-----
108 clean.data$d <- clean.data$donor
109 clean.data$d <- gsub('RRP2','MET-3B', clean.data$d)
110 clean.data$d <- gsub('RP2','MET-3A', clean.data$d)
111 clean.data$d <- gsub('RRP','MET-2B', clean.data$d)
112 clean.data$d <- gsub('RP','MET-2A', clean.data$d)
113 clean.data$r <- clean.data$run
114 clean.data$r <- gsub('30',1,clean.data$r)
115 clean.data$r <- gsub('33',3,clean.data$r)
116 clean.data$r <- gsub('31',2,clean.data$r)
117 clean.data$r <- gsub('32',4,clean.data$r)
118 #Performing PCA-----
119 pca.data <- perform_pca(clean.data, sampleID='sampleID', variable='variable',
120                          value='value',
121                          form.y=c('sampleID','run','vessel','donor','day','r','d'),
122                          form.x='variable', scale='UV', center=TRUE, convert.NA=TRUE
123                          )
124 #Exporting pca data for Kathleen-----
125 if(FALSE){
126   pca.names <- names(pca.data)
127   for(i in 1:length(pca.names)){
128     write.csv(file=paste('../Results/', pca.names[i], '-triplicates.csv', sep=''),
129              pca.data[[i]])
130   }
131   pca.summary <- as.matrix(summary(pca.data))
132   write.csv(file='../Results/summary-triplicates.csv', pca.summary[[6]])
133 }
134 #performing diagnostics on PCA-----
135 diagnostic.results = scree_plot(pca.data)
136
137 #unpacking the diagnostic results
138 diagnostic.number = diagnostic.results[['number']]
139 diagnostic.explained = diagnostic.results[['explained']]
140 diagnostic.plot = diagnostic.results[['plot']]
141 diagnostic.plot = diagnostic.plot + ggtitle('Scree Plot (UVN)')
142
143 #printing out diagnostic results
144 print(diagnostic.number)
145 print(diagnostic.explained)
146 print(diagnostic.plot)
147
148 #plotting PCA-processed data-----
149 #Score plot 2D
150 # Note: Arguments for score plot function are
151 #       (pca_data, qualifier_data, PCs=1:4, colour=NULL, shape=NULL, size=NULL,
152 #       label=NULL, title=NULL, legend.ori='horizontal', black-white=FALSE)
153
154 score_plot = score_plot_2d(pca.data, PCs=1:4,
155                            colour='r', shape='d', label='day',
156                            legend.ori='horizontal',
157                            black-white=FALSE)
158
159 p.score <- score_plot[['p']]
160 p.score <- p.score +
161   scale_colour_manual(values=c('forestgreen','darkorange2'),name='Run') + #set 3
162   #scale_colour_manual(values=c('coral2','darkorange2',
163   #                             'cornflowerblue','forestgreen'),name='Run') +
164   scale_shape_discrete(name='Donor')
165 collected.data <- score_plot[['collected.data']]

```

```

166
167 ggsave(' ../Plots/pca_plots/replicates3.png', p.score, height=15, width=25, units='
      cm')
168
169 #Loading plot 2D—————→
170
171 # Note: Arguments for loading plot function are (pca_data, PCs=1:4, label=NULL,
      title=NULL,
172 # view_mode='compressed')
173
174 #If you want all compound labels shown, enter all_paths in label argument
175 all_paths = rownames(pca.data$rotation)
176
177 #If you want all and only compound within the subset of observations processed in
      the PCA function, can create a temporary group called 'a' to enter in for the
      label argument
178 #a = pathway_list$Carb[pathway_list$Organic %in% rownames(pca_processed_data$
      rotation)]
179
180 loading_plot = loading_plot_2d(pca.data, PCs=1:4, label=all_paths,
181 title='2D Loading Plot-\nUV', view_mode = '
      compressed')
182
183 p.load <- loading_plot[[ 'p' ]]
184 load.data <- loading_plot[[ 'collected.data' ]]
185
186 print(p.load)
187 ggsave(' ../Results/set1-loading.png', p.load, height=30, width=50, units='cm')

```

profile diversity analysis.R

```

1 #####
2 #Profile diversity of Biomass samples
3
4 #sourcing required functions
5 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
6 source('function-data cleanup.R')
7 source('function-descriptive statistics.R')
8
9 #loading required packages
10 require(RSQLite)
11 require(ggplot2)
12 require(gridExtra)
13 require(RODBC)
14
15 #setting working directory
16 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass Samples/R Scripts/Analysis')
17
18 #Input data—————
19 #Choose driver
20 drv <- dbDriver('SQLite')
21
22 #Create connection (for SQLite, this is a file)
23 con <- dbConnect(drv, '../Data source files/data.db')
24
25 cmpd.info <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
26                               d.variable, d.value, c.Type
27                               FROM Cmpd_info c INNER JOIN Data d
28                               ON c.Compound=d.variable
29                               INNER JOIN Qualifier q
30                               ON d.file=q.file
31                               INNER JOIN Rank r
32                               ON r.Compound=d.variable
33                               WHERE r.Confidence = 1
34                               ORDER BY q.sampleID")
35
36 dbDisconnect(con)
37

```

```

38 # Create connection and selecting database
39 con <- odbcConnect('MySQLDSN')
40
41 input <- sqlQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
42                          d.variable, d.value, t.tax_class
43                          FROM biomass.data d INNER JOIN biomass.qualifier q
44                          ON d.file=q.file
45                          INNER JOIN general.tax t
46                          ON d.variable=t.syn
47                          INNER JOIN biomass.rank r
48                          ON r.compound=d.variable
49                          WHERE r.confidence = 1
50                          ORDER by q.sampleID;")
51
52 odbcClose(con)
53 #data cleanup-----
54 #Remove DSS, contaminants and compounds of zero concentration
55 clean.data <- data_cleanup(input, 'sampleID', 'variable', 'value',
56                             form.y=c('sampleID', 'run', 'vessel', 'donor', 'day', 'Type')
57                                     ,
58                                     form.x='variable', contaminants=NULL, export=FALSE,
59                                     convert.NA=TRUE)
60 #Calculating mean values of triplicates in set2 and input-----
61 if(any(clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3'))){
62   H.rep <- clean.data[clean.data$sampleID %in% c('H1', 'H2', 'H3'),]
63   R.rep <- clean.data[clean.data$sampleID %in% c('R1', 'R2', 'R3'),]
64
65   H <- ddply(H.rep, 'variable', mutate, avg=mean(value))
66   H$sampleID <- 'H'
67   H <- unique(H[, -7])
68   colnames(H)[8] <- 'value'
69
70   R <- ddply(R.rep, 'variable', mutate, avg=mean(value))
71   R$sampleID <- 'R'
72   R <- unique(R[, -7])
73   colnames(R)[8] <- 'value'
74
75   #removing H1,H2,H3,R1,R2,R3
76   clean.data <- clean.data[!clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3')
77                             ,]
78   #adding mean H and R
79   clean.data <- rbind(clean.data, H, R)
80   print('here')
81 }
82 #Adding extra qualifier to id RP-based cultures and RP2-based cultures-----
83 clean.data$culture.type <- 'RP'
84 clean.data$culture.type[clean.data$run %in% c(31,32)] <- 'RP2'
85 clean.data$culture.type[clean.data$run == 'Media control'] <- 'media'
86
87 #adding columns of coded run and donor labels-----
88 clean.data$d <- clean.data$donor
89 clean.data$d <- gsub('RRP2', 'MET-3B', clean.data$d)
90 clean.data$d <- gsub('RP2', 'MET-3A', clean.data$d)
91 clean.data$d <- gsub('RRP', 'MET-2B', clean.data$d)
92 clean.data$d <- gsub('RP', 'MET-2A', clean.data$d)
93 clean.data$r <- clean.data$run
94 clean.data$r <- gsub('30', 1, clean.data$r)
95 clean.data$r <- gsub('33', 3, clean.data$r)
96 clean.data$r <- gsub('31', 2, clean.data$r)
97 clean.data$r <- gsub('32', 4, clean.data$r)
98
99 #getting stats on each culture type
100 stat <- ddply(clean.data, c('variable', 'culture.type'), descriptive.stat,
101               data.type='long', variable='variable')
102
103 #total compound concentrations-----

```



```

104 rp <- clean.data[clean.data$donor=='RP',]
105 rrp <- clean.data[clean.data$donor=='RRP',]
106 rp2 <- clean.data[clean.data$donor=='RP2',]
107 rrp2 <- clean.data[clean.data$donor=='RRP2',]
108 total <- data.frame(rp30=sum(rp$value[ rp$run=='30' ]),
109                      rp33=sum(rp$value[ rp$run=='33' ]),
110                      rrp30=sum(rrp$value[ rrp$run=='30' ]),
111                      rrp33=sum(rrp$value[ rrp$run=='33' ]),
112                      rp2.31=sum(rp2$value[ rp2$run=='31' ]),
113                      rp2.32=sum(rp2$value[ rp2$run=='32' ]),
114                      rrp2.31=sum(rrp2$value[ rrp2$run=='31' ]),
115                      rrp2.32=sum(rrp2$value[ rrp2$run=='32' ]))
116 write.csv(file='../Results/total profile concentrations3.csv', total)
117
118 #Compound info - proportion of compound types-----
119 cmpd.info2 <- unique(clean.data[,c('variable','donor','tax_class','d')])
120 cmpd.table <- table(cmpd.info2$tax_class)/5
121 perc.table <- cmpd.table/sum(cmpd.table)
122
123 #proportion of compound types weighted to concentration-----
124 custom_perc <- function(data) {
125   output <- dplyr::ddply(data, 'sampleID', mutate, profile.total=sum(value, na.rm=TRUE))
126   out <- dplyr::ddply(output, c('sampleID','tax_class'), mutate,
127                          tax_class.total=sum(value, na.rm=TRUE))
128   out <- dplyr::ddply(out, 'tax_class', mutate, tax_class.avg= mean(unique(tax_class.total
129                                     )))
130   out$perc <- round(out$tax_class.total/out$profile.total*100,3)
131   out <- dplyr::ddply(out, 'tax_class', mutate, perc.avg=mean(unique(perc)))
132   out$perc.avg <- round(out$perc.avg, 3)
133   return(out)
134 }
135 cmpd.info3 <- dplyr::ddply(stat[stat$sampleID != 'Y',], 'donor', custom_perc)
136
137 cmpd.info4 <- unique(cmpd.info3[,c('donor','d','tax_class','perc.avg')])
138 cmpd.info4$donor <- factor(cmpd.info4$donor, levels=c('RP','RRP','RP2','RRP2'))
139
140 #cast for export
141 cmpd.info5 <- dcast(cmpd.info4, tax_class~donor, value.var='perc.avg')
142 cmpd.info5 <- cmpd.info5[order(cmpd.info5$RP, decreasing=TRUE),]
143
144 #write.csv(cmpd.info5, file='../Results/profile diversity3.csv')
145
146 #looking at compounds that differ between culture types-----
147 var.carb <- as.character(unique(cmpd.info3$variable[
148   cmpd.info3$tax_class == 'Carboxylic Acids and Derivatives']))
149
150 carb.acid <- unique(cmpd.info3[cmpd.info3$variable %in% var.carb,
151                               c('donor','run','culture.type','variable','tax_class',
152                                 'avg','se','profile.total')])
153
154 var.fat <- as.character(unique(cmpd.info3$variable[
155   cmpd.info3$tax_class == 'Fatty Acids and Conjugates']))
156
157 fatty.acid <- unique(cmpd.info3[cmpd.info3$variable %in% var.fat,
158                               c('donor','run','culture.type','variable','tax_class',
159                                 'avg','se','profile.total')])
160
161 #percentage of individual compounds
162 carb.perc <- carb.acid[order(carb.acid$variable),]
163
164 #first getting mean profile totals for each culture type
165 carb.perc <- dplyr::ddply(carb.perc, c('variable','culture.type'), mutate,
166                          avg.prof.tot = mean(profile.total))
167
168 carb.perc <- unique(carb.perc[,colnames(carb.perc) != 'profile.total'])
169 carb.perc$cmpd.perc <- carb.perc$avg/carb.perc$avg.prof.tot*100
170

```

```

171 fat.perc <- fatty.acid[order(fatty.acid$variable),]
172
173 #first getting mean profile totals for each culture type
174 fat.perc <- ddply(fat.perc, c('variable', 'culture.type'), mutate,
175                   avg.prof.tot = mean(profile.total))
176
177 fat.perc <- unique(fat.perc[, colnames(fat.perc) != 'profile.total'])
178 fat.perc$cmpd.perc <- fat.perc$avg/fat.perc$avg.prof.tot*100
179
180 #visualizing proportions-----
181 plot.data <- unique(cmpd.info3[, c('d', 'tax_class', 'tax_class.avg', 'perc.avg')])
182 #plot.data$donor <- factor(plot.data$donor, levels=c('RP', 'RRP', 'RP2', 'RRP2'))
183 plot.data <- plot.data[order(plot.data$tax_class.avg),]
184 plot.data$tax_class <- factor(plot.data$tax_class, levels=unique(plot.data$tax_
185                             class), ordered=TRUE)
186
187 fill.val <- rev(c('coral2', 'dodgerblue4', 'seagreen', 'yellow2', 'hotpink2', 'brown3',
188                  'cornflowerblue', 'darkorange2', 'seagreen2', 'purple3', 'tan4',
189                  'bisque3', 'slategray'))
190
191 p <- ggplot(plot.data, aes(x=d)) +
192   geom_bar(stat='identity', aes(y=tax_class.avg, fill=tax_class)) +
193   scale_fill_manual(values=fill.val,
194                     guide=guide_legend(reverse=TRUE,
195                                         title='Compound Class')) +
196   xlab('Community Culture') +
197   ylab('Class Structure of Metabolite Profiles\n(weighted to concentration)') +
198   theme_bw(9) +
199   theme(legend.key.size=unit(0.35, 'cm'))
200
201 ggsave('../Plots/profile_diversity3.png.png', p, height=10, width=15, units='cm')

```

donor comparison.R

```

1 #####
2 #Donor comparison
3
4 #loading required packages
5 require(ggplot2)
6 require(reshape2)
7 require(plyr)
8 require(RSQLite)
9
10 #Sourcing required functions
11 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions/')
12 source('function-data_cleanup.R')
13 source('function-descriptive_statistics.R')
14
15 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass samples/R Scripts/Analysis/')
16
17 #Input data-----
18 #Choose driver
19 drv <- dbDriver('SQLite')
20
21 #Create connection (for SQLite, this is a file)
22 con <- dbConnect(drv, '../Data source files/data.db')
23
24 input <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
25                          d.variable, d.value
26                          FROM Data d INNER JOIN Qualifier q
27                          ON d.file=q.file
28                          INNER JOIN Rank r
29                          ON d.variable = r.Compound
30                          WHERE r.Confidence = 1 AND q.run != 'Media control'
31                          ORDER BY q.sampleID;")
32
33 dbDisconnect(con)
34

```

```

35 #Data clean up-----
36 #Remove DSS, contaminants and compounds of zero concentration
37 clean.data <- data_cleanup(input, 'sampleID', 'variable', 'value',
38                             form.y=c('sampleID', 'run', 'vessel', 'donor', 'day'),
39                             form.x='variable', contaminants=NULL, export=FALSE,
40                             convert.NA=TRUE)
41
42 #Calculating mean values of triplicates in set2 and input-----
43 if(any(clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3'))){
44   H.rep <- clean.data[clean.data$sampleID %in% c('H1', 'H2', 'H3'),]
45   R.rep <- clean.data[clean.data$sampleID %in% c('R1', 'R2', 'R3'),]
46
47   H <- ddply(H.rep, 'variable', mutate, avg=mean(value))
48   H$sampleID <- 'H'
49   H <- unique(H[, -7])
50   colnames(H)[7] <- 'value'
51
52   R <- ddply(R.rep, 'variable', mutate, avg=mean(value))
53   R$sampleID <- 'R'
54   R <- unique(R[, -7])
55   colnames(R)[7] <- 'value'
56
57   #removing H1, H2, H3, R1, R2, R3
58   clean.data <- clean.data[!clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3')
59   ],]
60   #adding mean H and R
61   clean.data <- rbind(clean.data, H, R)
62   print('here')
63 }
64 #adding columns of coded run and donor labels-----
65 clean.data$d <- clean.data$donor
66 clean.data$d <- gsub('RRP2', 'MET-3B', clean.data$d)
67 clean.data$d <- gsub('RP2', 'MET-3A', clean.data$d)
68 clean.data$d <- gsub('RRP', 'MET-2B', clean.data$d)
69 clean.data$d <- gsub('RP', 'MET-2A', clean.data$d)
70 clean.data$r <- clean.data$run
71 clean.data$r <- gsub('30', 1, clean.data$r)
72 clean.data$r <- gsub('33', 3, clean.data$r)
73 clean.data$r <- gsub('31', 2, clean.data$r)
74 clean.data$r <- gsub('32', 4, clean.data$r)
75
76 #casting data to give non-profiled compounds concentration of zero
77 data.cast <- dcast(clean.data, sampleID+run+vessel+donor+day+d+r~variable,
78                    value.var='value')
79
80 data.cast[is.na(data.cast)] <- 0
81
82 #melting back to work with
83 data.melt <- melt(data.cast, id.vars=c('sampleID', 'run', 'vessel', 'donor', 'day', 'd',
84   'r'))
85
86 #adding ID for donor-run combo
87 data.melt$donor.run <- paste(data.melt$donor, data.melt$run, sep='.')
88
89 #Descriptive statistics-----
90 stat <- descriptive.stat(data.melt, data.type='long', variable='variable')
91
92 #Plot of RP vs RP2-----
93 plot.data <- stat
94 plot.data$x.code <- NULL
95 plot.data$x.code[plot.data$donor %in% c('RP', 'RRP')] <- 1
96 plot.data$x.code[plot.data$donor %in% c('RP2', 'RRP2')] <- 2
97 plot.data$x <- plot.data$x
98 plot.data$x[plot.data$run=='30'] <- plot.data$x[plot.data$run=='30'] - 0.18
99 plot.data$x[plot.data$run=='33'] <- plot.data$x[plot.data$run=='33'] + 0.18
100 plot.data$x[plot.data$run=='31'] <- plot.data$x[plot.data$run=='31'] - 0.18
101 plot.data$x[plot.data$run=='32'] <- plot.data$x[plot.data$run=='32'] + 0.18

```

```

101
102 #plot.data$run <- factor(plot.data$run, levels=c('30','33','31','32'))
103 #plot.data$donor <- factor(plot.data$donor, levels=c('RP','RRP','RP2','RRP2'))
104 culture <- plot.data[plot.data$donor %in% c('RP','RP2'),]
105 reciprocal <- plot.data[plot.data$donor %in% c('RRP','RRP2'),]
106
107
108 p <- ggplot(plot.data, aes(x=x, y=value, colour=r, linetype=d)) +
109   geom_point(aes(ymax=ymax, ymin=ymin)) +
110   stat_boxplot(data=culture, alpha=0.6) +
111   stat_boxplot(data=reciprocal, alpha=0.6) +
112   scale_x_continuous(breaks=c(1,2), label=c('MET-2','MET-3')) +
113   scale_linetype_manual(values=c(1,3,1,3), name='Defined culture') +
114   theme_bw(12) +
115   facet_wrap(~variable, scale='free_y', ncol=6) +
116   scale_colour_manual(values=c('lightcoral','chartreuse4','cyan3','darkorange'),
117                        name='Run') +
118   xlab('Defined Culture') +
119   ylab('Concentration (mM)') +
120   theme(legend.position='bottom')
121
122
123 ggsave('.. /Plots/donor differences4.png', p, height=25, width=23, units='cm')

```

time course analysis.R

```

1 #####
2 #Time-course analysis of Biomass samples
3
4 #sourcing required functions
5 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions')
6 source('function-data cleanup.R')
7 source('.. /Functions/function-extract legend.R')
8 source('.. /Functions/function-descriptive statistics.R')
9
10 #loading required packages
11 require(RSQLite)
12 require(ggplot2)
13 require(gridExtra)
14 require(plyr)
15 require(grid)
16
17
18 #setting working directory
19 setwd('C:/Users/Sandi/Dropbox/Projects/Biomass Samples/R Scripts/Analysis')
20
21 #Input data-----
22 #Choose driver
23 drv <- dbDriver('SQLite')
24
25 #Create connection (for SQLite, this is a file)
26 con <- dbConnect(drv, '.. /Data source files/data.db')
27
28 input <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
29                           d.variable, d.value
30                           FROM Data d INNER JOIN Qualifier q
31                           ON d.file=q.file
32                           INNER JOIN Rank r
33                           ON r.Compound=d.variable
34                           WHERE r.Confidence = 1
35                           ORDER BY q.sampleID;")
36
37 input2 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
38                            d.variable, d.value
39                            FROM Data d INNER JOIN Qualifier q
40                            ON d.file=q.file
41                            INNER JOIN Rank r
42                            ON r.Compound=d.variable

```

```

43         WHERE r.Confidence = 1
44         AND q.run != 'Media control'
45         ORDER BY q.sampleID;" )
46
47 set1 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
48                          d.variable, d.value
49                          FROM Data d INNER JOIN Qualifier q
50                          ON d.file=q.file
51                          INNER JOIN Rank r
52                          ON r.Compound=d.variable
53                          WHERE r.Confidence = 1
54                          AND q.run IN (30,33, 'Media control')
55                          ORDER BY q.sampleID;" )
56
57 set2 <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
58                          d.variable, d.value
59                          FROM Data d INNER JOIN Qualifier q
60                          ON d.file=q.file
61                          INNER JOIN Rank r
62                          ON r.Compound=d.variable
63                          WHERE r.Confidence = 1
64                          AND q.run IN (31,32, 'Media control')
65                          ORDER BY q.sampleID;" )
66
67 cmpd.info <- dbGetQuery(con, "SELECT q.sampleID, q.run, q.vessel, q.donor, q.day,
68                               d.variable, d.value, c.Type
69                               FROM Cmpd_info c INNER JOIN Data d
70                               ON c.Compound=d.variable
71                               INNER JOIN Qualifier q
72                               ON d.file=q.file
73                               ORDER BY q.sampleID" )
74
75 dbDisconnect(con)
76
77 #data cleanup-----
78 #Remove DSS, contaminants and compounds of zero concentration
79 clean.data <- data_cleanup(set1, 'sampleID', 'variable', 'value',
80                             form.y=c('sampleID', 'run', 'vessel', 'donor', 'day', 'Type')
81
82                             ,
83                             form.x='variable', contaminants=NULL, export=FALSE,
84                             convert.NA=TRUE)
85
86 #Calculating mean values of triplicates in set2 and input-----
87 if(any(clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3')) ) {
88   H.rep <- clean.data[clean.data$sampleID %in% c('H1', 'H2', 'H3'),]
89   R.rep <- clean.data[clean.data$sampleID %in% c('R1', 'R2', 'R3'),]
90
91   H <- ddply(H.rep, 'variable', mutate, avg=mean(value))
92   H$sampleID <- 'H'
93   H <- unique(H[, -7])
94   colnames(H)[7] <- 'value'
95
96   R <- ddply(R.rep, 'variable', mutate, avg=mean(value))
97   R$sampleID <- 'R'
98   R <- unique(R[, -7])
99   colnames(R)[7] <- 'value'
100
101 #removing H1, H2, H3, R1, R2, R3
102 clean.data <- clean.data[!clean.data$sampleID %in% c('H1', 'H2', 'H3', 'R1', 'R2', 'R3'),]
103 #adding mean H and R
104 clean.data <- rbind(clean.data, H, R)
105
106 print('here')
107 }
108 #adding columns of coded run and donor labels-----

```

```

109 clean.data$d <- clean.data$donor
110 clean.data$d <- gsub('RRP2','MET-3B', clean.data$d)
111 clean.data$d <- gsub('RP2', 'MET-3A', clean.data$d)
112 clean.data$d <- gsub('RRP', 'MET-2B', clean.data$d)
113 clean.data$d <- gsub('RP', 'MET-2A', clean.data$d)
114 clean.data$r <- clean.data$run
115 clean.data$r <- gsub('30',1,clean.data$r)
116 clean.data$r <- gsub('33',3,clean.data$r)
117 clean.data$r <- gsub('31',2,clean.data$r)
118 clean.data$r <- gsub('32',4,clean.data$r)
119
120 #casting data to give non-profiled compounds concentration of zero
121 data.cast <- dcast(clean.data, sampleID+run+vessel+donor+day+d+r~variable,
122                   value.var='value')
123
124 data.cast[is.na(data.cast)] <- 0
125
126 #melting back to work with
127 data.melt <- melt(data.cast, id.vars=c('sampleID','run','vessel','donor','day','d',
128   'r'))
129
130 #adding ID for donor-run combo
131 data.melt$donor.run <- paste(data.melt$donor, data.melt$run, sep='.')
132
133 #ID for run replicates
134 data.melt$run.rep <- data.melt$donor.run
135 data.melt$run.rep[data.melt$run.rep %in% c('RP.30','RP.33','RRP.30','RRP.33')] <-
136   '30.33'
137 data.melt$run.rep[data.melt$run.rep %in% c('RP2.31','RP2.32','RRP2.31','RRP2.32')]
138   <- '31.32'
139
140 #converting day to numeric
141 data.melt$day <- as.numeric(data.melt$day)
142
143 #Plotting over time-bar


---


144 plot.list <- list()
145 for(i in unique(data.melt$variable)) {
146   plot.data <- data.melt[data.melt$variable == i,]
147   plot.data$run <- factor(plot.data$run, levels=c('Media control','30','33','31','32'))
148   plot.list[[i]] <- ggplot(plot.data, aes(x=donor, y=value, fill=day, shape=run)) +
149     geom_bar(position='dodge', stat='identity') +
150     facet_wrap(~ run, scales='free_x', nrow=1) +
151     theme_bw(12) +
152     ggtitle(i) +
153     theme(legend.position='none',
154           title=element_text(size=8),
155           axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
156 }
157
158 l.plot <- ggplot(plot.data, aes(x=donor, y=value, fill=day, shape=run)) +
159   geom_bar(position='dodge', stat='identity') +
160   facet_wrap(~ run, scales='free_x', nrow=1) +
161   theme_bw(14)
162 legend <- g_legend(l.plot)
163
164 png(filename='../Plots/Time-bar plots.png', width=50, height=60, units='cm', res
165   =240)
166 do.call(grid.arrange, plot.list)
167 dev.off()
168
169 #Looking for linearity


---


170 custom_linear <- function(data, group) {
171   aov.lm <- dplyr::dplyr(data, group, aov, formula=value~day)
172   out <- lapply(aov.lm, anova)

```

```

170   return(out)
171 }
172
173 lm.run <- ddply(data.melt, 'variable', custom_linear, 'run')
174 s.lm.run <- lm.run[lm.run[, 'Pr(>F)'] <= 0.1 & !is.na(lm.run[, 'Pr(>F)']) ,]
175
176 lm.dr <- ddply(data.melt, 'variable', custom_linear, 'donor.run')
177 s.lm.dr <- lm.dr[lm.dr[, 'Pr(>F)'] <= 0.1 & !is.na(lm.dr[, 'Pr(>F)']) ,]
178
179 #Plotting over time-scatter
180
181 p.scatter <- ggplot(data.melt, aes(x=day, y=value, shape=donor, colour=run)) +
182   geom_point() +
183   geom_path() +
184   facet_wrap(~ variable, scales='free') +
185   scale_colour_manual(values=c('chartreuse4', 'darkorange', 'blue')) +
186   theme_bw(12) +
187   theme(title=element_text(size=8),
188         axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
189
190 ggsave(file='../Plots/Time course/Run 31 and 32/All compounds2.png', p.scatter,
191        height=25, width=30, units='cm')
192
193 #Compounds showing linearity as defined by runs-----
194 lm.run <- data.melt[data.melt$variable %in% unique(s.lm.run$variable),]
195
196 p.lm.run <- ggplot(lm.run, aes(x=day, y=value, shape=donor, colour=run)) +
197   geom_point() +
198   geom_path() +
199   facet_wrap(~ variable, scales='free') +
200   theme_bw(12) +
201   theme(title=element_text(size=8),
202         axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
203
204 ggsave(file='../Plots/Time course/Run 31 and 32/linear-run.png', p.lm.run,
205        height=15, width=25, units='cm')
206
207 #Compounds showing linearity as defined by donors within runs-----
208 lm.dr <- data.melt[data.melt$variable %in% unique(s.lm.dr$variable),]
209
210 p.lm.dr <- ggplot(lm.dr, aes(x=day, y=value, shape=donor, colour=run)) +
211   geom_point() +
212   geom_path() +
213   facet_wrap(~ variable, scales='free') +
214   theme_bw(12) +
215   theme(title=element_text(size=8),
216         axis.text.x=element_text(hjust=0, vjust=1, angle=-45))
217
218 ggsave(file='../Plots/Time course/Run 30 and 33/linear-donor within runs.png', p.lm
219        .dr,
220        height=20, width=30, units='cm')
221
222 #Reaching common point Run 30 and 33-----
223 comm.lm <- c('Alanine', 'Asparagine', 'Aspartate', 'Betaine', 'Ethanol', 'Fructose',
224             'Glycerol', 'Glycolate', 'Hydroxyacetone', 'Propionate', 'Thymine',
225             'Uracil', 'Valerate')
226
227 comm.notlm <- c('Glucose', 'Phenylalanine', 'Succinate', 'Threonine', 'Tyrosine',
228               'Xylose')
229
230 one.comm <- c('Acetate', 'Butyrate', 'Choline', 'Fumarate', 'Glycine', 'Isobutyrate',
231             'Isoleucine', 'Leucine', 'Lysine', 'Methanol', 'Methionine', 'Methylamine',
232             'Proline', 'Serine', 'Trimethylamine', 'Valine')
233
234 cmpd.list <- list(comm.lm, comm.notlm, one.comm)

```

```

235 p.list <- list()
236 ncol <- c(4,3,4)
237 width <- c(17, 15, 20)
238 height <- c(12, 10, 20)
239
240 for(i in 1:length(cmpd.list)) {
241   plot.data <- data.melt[data.melt$variable %in% cmpd.list[[i]],]
242   plot.data <- plot.data[plot.data$run != 'Media control',]
243
244   p.list[[i]] <- ggplot(plot.data, aes(x=day, y=value, shape=d, colour=r)) +
245     geom_point() +
246     geom_path() +
247     theme_bw(11) +
248     facet_wrap(~variable, scales='free', ncol=ncol[i]) +
249     scale_colour_discrete(name='Run') +
250     scale_shape_discrete(name='Donor') +
251     scale_x_continuous(breaks=c(6,14,20)) +
252     xlab('Day') +
253     ylab('Concentration (mM)') +
254     theme(legend.direction='horizontal', legend.position='bottom',
255           legend.box='horizontal')
256 }
257
258 vp <- viewport(x=0.5, y=0.5, width=unit(17,'cm'), height=unit(18,'cm'))
259
260 png(filename='../Plots/Time course/Run 30 and 33/common point2.png',
261     width=17, height=25, units='cm', res=300, pointsize=11)
262 plot.new()
263 pushViewport(vp)
264 print(p.list[[1]], vp=viewport(x=0.5, y=0.75, width=unit(17,'cm'),
265                               height=unit(16,'cm'))))
266 print(p.list[[2]], vp=viewport(x=0.385, y=0.12, width=unit(13,'cm'),
267                               height=unit(10,'cm'))))
268 grid.text('A', x=0.03, y=1.155, gp=gpar(cex=1.5))
269 grid.text('B', x=0.03, y=0.36, gp=gpar(cex=1.5))
270 dev.off()
271
272 ggsave('../Plots/Time course/Run 30 and 33/common point-one run2.png', p.list[[3]],
273         width=18, height=18, units='cm')
274
275 #Reaching common point Run 31 and 32
276 comm.lm <- c('Alanine', 'Asparagine', 'Glutamate', 'Glycine', 'Proline')
277
278 comm.notlm <- c('Butyrate', 'Glucose', 'Galactose', 'Glutamate', 'Hydroxyacetone',
279               'Isobutyrate', 'Isoleucine', 'Isovalerate', 'Methionine',
280               'Phenylalanine', 'Propionate', 'Tyrosine', 'Uracel', 'Valerate',
281               'Valine', 'Xylose')
282
283 one.comm <- c('Acetate', 'Aspartate', 'Benzoate', 'Betaine', 'Choline', 'Ethanol',
284             'Leucine', 'Methanol', 'Phenylacetate', 'Pyroglutamate',
285             'Trimethylamine')
286
287 cmpd.list <- list(comm.lm, comm.notlm, one.comm)
288
289 p.list <- list()
290 ncol <- c(5,5,4)
291 width <- c(20, 17, 18)
292 height <- c(5, 16, 12)
293 for(i in 1:length(cmpd.list)) {
294   plot.data <- data.melt[data.melt$variable %in% cmpd.list[[i]],]
295   plot.data <- plot.data[plot.data$run != 'Media control',]
296
297   p.list[[i]] <- ggplot(plot.data, aes(x=day, y=value, shape=d, colour=r)) +
298     geom_point() +
299     geom_path() +
300     theme_bw(11) +
301     facet_wrap(~variable, scales='free', ncol=ncol[i]) +
302     scale_colour_manual(values=c('chartreuse4', 'darkorange'), name='Run') +

```



```

303     scale_shape_discrete(name='Donor') +
304     scale_x_continuous(breaks=c(6,14,20)) +
305     xlab('Day') +
306     ylab('Concentration (mM)') +
307     theme(legend.direction='horizontal', legend.position='bottom',
308           legend.box='horizontal')
309 }
310
311 vp <- viewport(x=0.5, y=0.5, width=unit(17,'cm'), height=unit(15,'cm'))
312
313 png(filename='../Plots/Time course/Run 31 and 32/common point2.png',
314      width=17, height=25, units='cm', res=300, pointsize=11)
315 plot.new()
316 pushViewport(vp)
317 print(p.list[[1]], vp=viewport(x=0.5, y=1.1, width=1, height=0.4))
318 print(p.list[[2]], vp=viewport(x=0.5, y=0.57, width=1, height=0.85))
319 grid.text('A', x=0.03, y=1.3, gp=gpar(cex=1.5))
320 grid.text('B', x=0.03, y=0.95, gp=gpar(cex=1.5))
321 dev.off()
322
323 ggsave('../Plots/Time course/Run 31 and 32/common point-one run2.png', p.list[[3]],
324        width=width[3], height=15, units='cm')

```

C.4 Metabolomic Analysis of Human Fecal Microbiota

Data source.R

```

1  #####
2  #This is a database of all donors and repopulate cultures
3
4  #loading required packages
5  require(RSQLite)
6  require(reshape2)
7
8  #Setting working directory
9  setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Data source
    files/')
10
11 #Sourcing required functions
12 source('../functions/function-summary preprocessing.R')
13
14 #Getting raw data files from each project
15 #
16
17
18 rp2.ua.raw <- summary.preprocess(file.name='C:/Users/Sandi/Dropbox/Projects/Uric
    acid Project/Results/Expt 1 and 2-UA NMR/R Scripts/Data source files/UA Expt 2
    summary2.csv')
19
20 rp2.ua <- rp2.ua.raw[['data']]
21 molwt.rp2.ua <- rp2.ua.raw[['mol.wt']]
22
23 rp2.ua.melt <- melt(rp2.ua, id.vars='file')
24
25 d9.raw <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Uric acid Project/
    Results/Expt 1 and 2-UA NMR/R Scripts/Data source files/UA Expt 1 summary2.csv'
    )
26
27 d9 <- d9.raw[['data']]
28 molwt.d9 <- d9.raw[['mol.wt']]
29
30 d9.melt <- melt(d9, id.vars='file')

```

```

31 #
32 #Repopulate samples
33 rp1.raw <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Repopulate samples/
  R scripts/Data source files/Concentrations-RepoopI_mM3.csv')
34
35 rp1 <- rp1.raw[['data']]
36 molwt.rp1 <- rp1.raw[['mol.wt']]
37
38 rp1.melt <- melt(rp1, id.vars='file')
39
40 rp2.raw <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Repopulate samples/
  R scripts/Data source files/Concentrations-RepoopII_mM2.csv')
41
42 rp2 <- rp2.raw[['data']]
43 molwt.rp2 <- rp2.raw[['mol.wt']]
44
45 rp2.melt <- melt(rp2, id.vars='file')
46 #
47
48 #Liquid Gold samples
49 lg.raw <- summary.preprocess('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts
  /V 3.0/Data source files/D5_D8 summary_mM3.csv')
50
51 lg <- lg.raw[['data']]
52 molwt.lg <- lg.raw[['mol.wt']]
53
54 lg.melt <- melt(lg, id.vars='file')
55 #
56
57 #putting all samples together into one data frame
58 data <- rbind(rp2.ua.melt, d9.melt, rp1.melt, rp2.melt, lg.melt)
59
60 #getting molecular weight of compounds
61 cmpd.molwt <- unique(as.data.frame(t(cbind(molwt.rp2.ua, molwt.d9,
  molwt.rp1, molwt.rp2, molwt.lg))))
62 cmpd.molwt$variable <- rownames(cmpd.molwt)
63
64 #
65
66 #Making Qualifiers
67 q.rp2.ua <- data.frame(file=rp2.ua$file,
68   sampleID=LETTERS[1:17],
69   donor=rep('RP2',17),
70   vessel=c(0, rep(1,6), rep(2, 10)),
71   DPI=rep(NA, 17),
72   phase=c(0, rep(1,3), rep(2,3), rep(3,7), rep(4,3)),
73   hour=c(1(-24)
74     ,0,24,32,48,72,80,0.25,0.5,1,2,4,24,32,48,72,80),
75   dose=c(0,rep(360,3),rep(720,3),rep(360,10)),
76   dose.type=c(NA, 'spike',rep('feed',15)),
77   expt=rep('UA Expt2',17))
78
79 q.rp2.ua$scan <- NA
80 q.rp2.ua$treatment <- c(rep('UA',14),rep('Vancomycin',3))
81 q.rp2.ua$run <- 5
82 q.rp2.ua$gen.treat <- c(rep('None', 14), rep('Vancomycin',3))
83 q.rp2.ua$state <- 'ss'
84
85 #making donor of class character
86 q.rp2.ua$donor <- as.character(q.rp2.ua$donor)

```

```

85
86 q.d9 <- data.frame( file=d9$file ,
87                     sampleID=LETTERS[18:25] ,
88                     donor=rep( 'C' ,8) ,
89                     vessel=rep( 3,8) ,
90                     DPI=c( 40,41,44,45,48,49,52,53) ,
91                     phase=rep(NA,8) ,
92                     hour=rep(NA,8) ,
93                     dose=c( rep(0,2) , rep(360,6)) ,
94                     dose.type=rep( 'feed' ,8) ,
95                     expt=rep( 'UA Expt1' ,8))
96
97 #adding columns to match other qualifier sets
98 q.d9$scan <- NA
99 q.d9$treatment <- c(rep( 'None' ,2) ,rep( 'UA' ,nrow(q.d9)-2))
100 q.d9$run <- 4
101 q.d9$gen.treat <- rep( 'None' , nrow(q.d9))
102 q.d9$state <- 'ss'
103
104 #making donor of class character
105 q.d9$donor <- as.character(q.d9$donor)
106
107 q.rp1 <- data.frame( file=rp1$file ,
108                     sampleID=paste( 'RP1.D' ,c(0,3,5,7,10) , sep='') ,
109                     DPI=c( 0,3,5,7,10) ,
110                     donor=rep( 'RP1' ,5) ,
111                     expt=rep( 'RP1' ,5))
112
113 #adding columns to match other qualifier sets
114 q.rp1$vessel <- NA
115 q.rp1$phase <- NA
116 q.rp1$dose <- NA
117 q.rp1$dose.type <- NA
118 q.rp1$scan <- NA
119 q.rp1$treatment <- rep( 'none' ,nrow(q.rp1))
120 q.rp1$run <- 6
121 q.rp1$hour <- NA
122 q.rp1$gen.treat <- rep( 'None' , nrow(q.rp1))
123 q.rp1$state <- c( 'uss' , rep( 'ss' ,4))
124
125 #making donor of class character
126 q.rp1$donor <- as.character(q.rp1$donor)
127
128 q.rp2 <- data.frame( file=rp2$file ,
129                     sampleID= paste( 'RP2.V' ,c(1,2) , sep='') ,
130                     DPI=rep(NA, 2) ,
131                     donor=rep( 'RP2' ,2) ,
132                     expt=rep( 'RP2' , 2))
133
134 #adding columns to match other qualifier sets
135 q.rp2$vessel <- NA
136 q.rp2$phase <- NA
137 q.rp2$dose <- NA
138 q.rp2$dose.type <- NA
139 q.rp2$scan <- NA
140 q.rp2$treatment <- rep( 'None' , nrow(q.rp2))
141 q.rp2$run <- 7
142 q.rp2$hour <- NA
143 q.rp2$gen.treat <- rep( 'None' , nrow(q.rp2))
144 q.rp2$state <- 'ss'
145
146 #making donor of class character
147 q.rp2$donor <- as.character(q.rp2$donor)
148
149 q.lg <- read.csv( 'C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Data
150                  source files/Liquid gold qualifier.csv' ,
151                  header=TRUE, check.names=FALSE, encoding='UTF-8')
152 q.lg$expt <- NA

```

```

152 q.lg$expt[q.lg$Donor == 'A'] <- 'Liquid gold_D5'
153 q.lg$expt[q.lg$Donor == 'B'] <- 'Liquid gold_D8'
154 q.lg <- q.lg[, colnames(q.lg) != 'Plot.Treat']
155 colnames(q.lg) <- c('file', 'scan', 'donor', 'run', 'treatment', 'DPI', 'state', 'vessel',
156                    'gen.treat', 'sampleID', 'expt')
157
158 q.lg$run[q.lg$run == 1] <- 8
159 q.lg$run[q.lg$run == 2] <- 1
160 q.lg$run[q.lg$run == 3] <- 2
161 q.lg$run[q.lg$run == 4] <- 3
162
163 #adding columns to match other qualifier sets
164 q.lg$phase <- NA
165 q.lg$dose <- NA
166 q.lg$hour <- NA
167 q.lg$dose.type <- NA
168
169
170 #making donor of class character
171 q.lg$donor <- as.character(q.lg$donor)
172
173 #putting all qualifiers together
174 qualifier <- rbind(q.rp2.ua, q.d9, q.rp1, q.rp2, q.lg)
175 qualifier$expt <- as.character(qualifier$expt)
176 qualifier$vessel <- as.character(qualifier$vessel)
177 qualifier$state <- as.character(qualifier$state)
178
179 #

```

```

180 #Getting compound annotations
181 rank.rp2.ua <- read.csv('C:/Users/Sandi/Dropbox/Projects/Uric acid Project/Results/
    Expt 1 and 2-UA NMR/R Scripts/Data source files/Compound annotations-RepoopII-
    UA.csv',
182                      encoding='UTF-8')
183 rank.rp2.ua <- rank.rp2.ua[, 1:7]
184 rank.rp2.ua$expt <- 'UA Expt2'
185
186 rank.d9 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Uric acid Project/Results/Expt
    1 and 2-UA NMR/R Scripts/Data source files/Compound annotations-Donor 9.csv',
187                  encoding='UTF-8')
188 rank.d9 <- rank.d9[, 1:7]
189 rank.d9$expt <- 'UA Expt1'
190
191 rank.rp1 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Repoopulate samples/R scripts
    /Data source files/Compound annotations-RepoopI.csv',
192                  encoding='UTF-8')
193 rank.rp1 <- rank.rp1[, 1:7]
194 rank.rp1$expt <- 'RP1'
195
196 rank.rp2 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Repoopulate samples/R scripts
    /Data source files/Compound annotations-RepoopII.csv',
197                  encoding='UTF-8')
198 rank.rp2 <- rank.rp2[, 1:7]
199 rank.rp2$expt <- 'RP2'
200
201 rank.d5 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 2.0/
    data source files/Export-liquid gold confidence-D5.csv',
202                  header=TRUE, encoding='UTF-8')
203 rank.d5 <- rank.d5[, 1:7]
204 colnames(rank.d5) <- c('Compound', 'Clusters', 'Convolution', 'Consistency', '
    Confidence', 'X.. of. Profiles', 'X.. of. Profiles')
205 rank.d5$expt <- 'Liquid gold_D5'
206
207 rank.d8 <- read.csv('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 2.0/
    data source files/Export-liquid gold confidence-D8.csv',
208                  header=TRUE, encoding='UTF-8')
209 rank.d8 <- rank.d8[, 1:7]

```

```

210 colnames(rank.d8) <- c('Compound', 'Clusters', 'Convolution', 'Consistency', '
    Confidence', 'X..of.Profiles', 'X..of.Profiles')
211 rank.d8$expt <- 'Liquid gold_D8'
212
213
214 #putting all rankings together
215 rank <- rbind(rank.rp2.ua, rank.d9, rank.rp1, rank.rp2, rank.d5, rank.d8)
216 #

```

```

217 #metabolite pathway
218 path.data <- read.csv('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/
    Data source files/Compound pathway annotation.csv',
219     header=FALSE)
220 path.data <- path.data[1:153,]
221
222 #Getting compound types
223 cmpd.type <- data.frame(Compound=t(path.data[1,2:ncol(path.data)]),
224     Type=t(path.data[2,2:ncol(path.data)]))
225 colnames(cmpd.type) <- c('Compound', 'Type')
226
227 #Categorizing pathways
228 main.path <- data.frame(main=c(rep('Carbohydrates', 12),
229     rep('Catecholamines', 5),
230     rep('Vitamins', 8),
231     rep('Antibiotics', 3),
232     rep('Lipids', 9),
233     rep('Amino acids', 28),
234     rep('Purine, Pyrimidine', 4),
235     rep('Fatty acids', 6),
236     rep('Organic acids', 9),
237     rep('Proteins', 5),
238     rep('Inorganic Compounds', 4),
239     rep('Steroid Hormones', 4),
240     rep('Central Metabolic Pathways', 8),
241     rep('Aromatic Compounds', 8),
242     rep('Neurotransmitters', 3),
243     rep('Cytochromes', 3),
244     rep('Other', 31)),
245     path=path.data[4:nrow(path.data),1])
246
247 # removing pathway headings from excel version
248 main <- main.path[-c(1,13,18,26,29,38,66,70,
249     76,85,90,94,98,106,114,117,120),]
250
251 #Getting pathway assignment for each compound
252 col.name <- as.character(as.matrix(path.data[1, 2:ncol(path.data)]))
253 row.name <- as.character(as.matrix(path.data[4:nrow(path.data),1]))
254
255 path <- path.data[4:nrow(path.data), 2:ncol(path.data)]
256 colnames(path) <- col.name
257 rownames(path) <- row.name
258 path <- path[-c(1,13,18,26,29,38,66,70,
259     76,85,90,94,98,106,114,117,120),]
260
261 pathway <- c()
262 for(i in 1:ncol(path)) {
263     temp <- path[,i]
264     yes <- rownames(temp)[temp[1] == 'Y']
265     out <- data.frame(Compound=rep(colnames(temp), length(yes)), Pathway=yes)
266     pathway <- rbind(pathway, out)
267 }
268
269 # Create the connection (no file needed)
270
271 # Choose driver
272 drv = dbDriver('SQLite')
273

```

```

274 # Create connection (for SQLite, this is a file)
275 con = dbConnect(drv, 'data.db')
276
277 # Creating database tables
278 data.table <- dbWriteTable(con, 'Data', data, overwrite=TRUE)
279
280 qualifier.table <- dbWriteTable(con, 'Qualifier', qualifier, overwrite=TRUE)
281
282 q.rp2.ua.table <- dbWriteTable(con, 'Q_RP2-UA', q.rp2.ua, overwrite=TRUE)
283 q.d9.table <- dbWriteTable(con, 'Q-D9', q.d9, overwrite=TRUE)
284 q.rp1.table <- dbWriteTable(con, 'Q-RP1', q.rp1, overwrite=TRUE)
285 q.rp2.table <- dbWriteTable(con, 'Q-RP2', q.rp2, overwrite=TRUE)
286 q.lg.table <- dbWriteTable(con, 'Q-LG', q.lg, overwrite=TRUE)
287
288 rank.table <- dbWriteTable(con, 'Rank', rank, overwrite=TRUE)
289
290 molwt.table <- dbWriteTable(con, 'Mol_wt', compd.molwt, overwrite=TRUE)
291
292 compdtype.table <- dbWriteTable(con, 'Compd.type', compd.type, overwrite=TRUE)
293
294 main.table <- dbWriteTable(con, 'Main_pathway', main, overwrite=TRUE)
295
296 path.table <- dbWriteTable(con, 'Pathway', pathway, overwrite=TRUE)
297
298 tables <- dbGetQuery(con, "SELECT * FROM sqlite_master WHERE type='table';")
299 print(tables)
300
301 dbDisconnect(con)

```

PCA analysis.R

```

1 #####
2 #Modified for Manuscript A
3 #Analysis to compare samples across projects
4
5 #loading requiried packages
6 require(RSQLite)
7 require(ggplot2)
8
9 #sourcing required functions
10 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions/')
11 source('function-pca.R')
12 source('function-plotting 2D score.R')
13 source('function-plotting 2D loading.R')
14 source('function-pca diagnostics.R')
15 source('function-data cleanup.R')
16 source('function-scaling.R')
17
18 #setting working directory
19 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript A'
20 )
21 #
22 #Input data
23 # Choose driver
24 drv <- dbDriver('SQLite')
25
26 # Create connection (for SQLite, this is a file)
27 con <- dbConnect(drv, '../Data source files/data.db')
28
29 input <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
30 q.hour, q.phase, q.dose, q.treatment, q.gen_treat, q.expt,
31 d.variable, d.value
32 FROM Data d INNER JOIN Qualifier q
33 ON d.file = q.file
34 INNER JOIN Rank r
35 ON d.variable = r.Compound AND q.expt = r.expt
36 INNER JOIN Compd.type c

```

```

37         ON r.Compound = c.Compound
38         WHERE r.Confidence = 1 AND q.gen_treat = 'None'
39         AND (q.run != 8 OR q.run IS NULL)
40         AND d.variable NOT IN ('Xanthine', 'Glycerol')
41         AND q.state = 'ss'
42         AND q.donor != 'RP1'
43         AND q.expt != 'RP2'
44         ORDER BY q.sampleID;")
45
46 dbDisconnect(con)
47 #-----
48 #converting factors to appropriate classes
49 input$run <- as.character(input$run)
50
51 #-----
52 #Performing cleanup function to remove DSS, contaminants and compounds of
53 #zero concentration
54
55 clean.data <- data_cleanup(input, 'sampleID', 'variable', 'value',
56                             form.y=c('sampleID', 'donor', 'run', 'vessel', 'DPI',
57                                         'hour', 'phase', 'dose', 'treatment', 'gen_treat',
58                                         'expt'),
59                             form.x='variable', contaminants=NULL, export=FALSE)
60
61 clean.data$dr <- paste(clean.data$donor, clean.data$run, sep='.')
62 clean.data$dr <- factor(clean.data$dr, levels=c('A.1', 'A.2', 'B.3', 'C.4', 'RP2.5'))
63 #-----
64 #performing the PCA
65 pca.data = perform_pca(clean.data, sampleID='sampleID',
66                         variable='variable', value='value',
67                         form.y=c('sampleID', 'donor', 'run', 'vessel', 'DPI',
68                                     'hour', 'phase', 'dose', 'treatment', 'gen_treat',
69                                     'expt', 'dr'),
70                         form.x='variable',
71                         scale='UV', center=TRUE, convert.NA=TRUE)
72 #-----
73 #performing diagnostics on PCA
74 diagnostic.results = scree_plot(pca.data)
75
76 #unpacking the diagnostic results
77 diagnostic.number = diagnostic.results[['number']]
78 diagnostic.explained = diagnostic.results[['explained']]
79 diagnostic.plot = diagnostic.results[['plot']]
80 diagnostic.plot = diagnostic.plot + ggtitle('Scree Plot (UVN)')
81
82 #printing out diagnostic results
83 print(diagnostic.number)
84 print(diagnostic.explained)
85 print(diagnostic.plot)
86 #-----
87 #plotting PCA-processed data
88 #Score plot 2D
89 # Note: Arguments for score plot function are
90 #       (pca_data, qualifier_data, PCs=1:4, colour=NULL, shape=NULL, size=NULL,
91 #       label=NULL, title=NULL, legend.ori='horizontal', black_white=FALSE)
92
93 score_plot = score_plot_2d(pca.data, PCs=1:4,
94                             shape='dr',
95                             legend.ori='horizontal',
96                             outline=FALSE, black_white=FALSE)
97
98 p.score <- score_plot[['p']]
99 collected.data <- score_plot[['collected.data']]
100
101 ellipse_custom <- function(d) {
102
103     #Calculating ellipse for 95% confidence
104     ellipseSD = cor(d$x.value, d$y.value)

```

```

105 ellipseDf = as.data.frame(ellipse(ellipseSD, scale=c(sd(d$x.value),
106                                                         sd(d$y.value))),
107                             centre=colMeans(d[,c('x.value', 'y.value')]),
108                             t=2))
109 return(ellipseDf)
110 }
111
112 e.path <- ddply(collected.data, 'view', ellipse_custom)
113
114 #Drawing an ellipse for 95% confidence
115 p.score = p.score +
116   # geom_path(data=e.path, aes(x=x, y=y), linetype=2, alpha=0.5) +
117   # scale_colour_manual(values=c('black', 'darkorange1', 'cornflowerblue',
118   #                               'black', 'chartreuse', 'blue', 'grey30'),
119   #
120   #                       guide_legend(title='Run'),
121   #                       labels=c('NA', '1', '2', '3', '4', '5', '6')) +
122   #scale_colour_manual(values=c('deeppink1', 'darkorange1', 'cyan4',
123   #                               'darkorchid4', 'chartreuse', 'blue', 'grey30'),
124   #                       guide_legend(title='Run'), labels=c(1:5)) +
125   scale_shape_manual(name='', labels=c('Donor A (Run 1)', 'Donor A (Run 2)',
126   #                                     'Donor B', 'Donor C', 'MET-2'),
127   #                   values=c(16,17,7,3,5)) +
128   guides(col=guide_legend(byrow=TRUE))
129
130 ggsave('.. /Plots/PCA plots/Profiles-score5.png', p.score, height=15, width=25,
131         units='cm')
132
133 tiff(file='C:/Users/Sandi/Dropbox/Projects/Liquid Gold/Manuscript/Manuscript A/
134       Journal of Proteome Research/Figure 2.tif',
135       width=7, height=10/2.54, units='in', res=300,)
136 print(p.score)
137 dev.off()
138 #
139
140 #Loading plot 2D
141 # Note: Arguments for loading plot function are (pca_data, PCs=1:4, label=NULL,
142 #         title=NULL,
143 #         view_mode='compressed')
144 #If you want all compound labels shown, enter all_paths in label argument
145 all_paths = rownames(pca.data$rotation)
146
147 source('.. /Functions/defined binning.R')
148 bin <- defined_binning()
149
150 #If you want all and only compound within the subset of observations processed in
151 #the PCA function, can create a temporary group called 'a' to enter in for the
152 #label argument
153 #a = pathway_list$Carb[pathway_list$Organic %in% rownames(pca_processed_data$
154 #rotation)]
155 loading_plot = loading_plot_2d(pca.data, PCs=1:4, label=all_paths,
156                                title=NULL, outline=FALSE, view_mode = 'compressed')
157
158 p.load <- loading_plot[['p']]
159 load.data <- loading_plot[['collected.data']]
160
161 load.data$bin <- NA
162 for(i in bin$variable) {
163   load.data$bin[load.data$Text == i] <- as.character(bin$bin[bin$variable == i])
164 }
165
166 p.load <- p.load +
167   geom_point(data=load.data, aes(x=xPC, y=yPC, shape=bin), size=1) +
168   geom_text(data=load.data, aes(x=xPC, y=yPC, label=Text),
169             hjust=-0.1, vjust=0, size=1, colour='grey50') +
170   scale_colour_manual(values=c('deeppink1', 'darkorange1', 'cyan4',

```



```

165         'darkorchid4', 'chartreuse'),
166         guide_legend(title='Bin')) +
167     scale_shape_manual(name='Bin', values=rev(c(16,17,7,3,5))) +
168     theme(legend.position='bottom')
169
170 ggsave(' ../Plots/PCA plots/Profiles-loading5.png', p.load,
171        height=12*1.5, width=18*2, units='cm')
172
173 #EPS
174 p <- ggplot(load.data, aes(xPC, yPC, shape=bin)) +
175     geom_text(aes(label=Text), size=4, hjust=-0.1, vjust=0.4, colour='grey50') +
176     geom_point(size=3) +
177     scale_shape_manual(name='Bin', values=rev(c(16,17,7,3,5))) +
178     facet_wrap(~ View, ncol=min(floor(length(1:4)/2), 3)) +
179     xlab('PC score') +
180     ylab('PC score') +
181     theme_bw(12) +
182     theme(axis.title.x = element_text(size=14, face='bold'),
183           axis.title.y = element_text(size=14, face='bold'),
184           #axis.text = element_text(size=12, colour='black'),
185           legend.title = element_text(size=14),
186           legend.text = element_text(size=14),
187           legend.key = element_rect(colour='white'),
188           legend.margin = unit(0, "cm"),
189           panel.border=element_rect(size=2, colour='black'),
190           strip.text.x = element_text(size=12, face='bold'),
191           legend.position='bottom')
192 #Drawing origin
193 p = p + geom_vline(xintercept=0, alpha=0.3)
194 p = p + geom_hline(yintercept=0, alpha=0.3)
195
196 #Calculating plot limits
197 xlimit = max(abs(load.data$xPC))*1.15
198 ylimit = max(abs(load.data$yPC))*1.1
199
200 #Setting limits
201 p = p + xlim(-xlimit, xlimit)
202 p = p + ylim(-ylimit, ylimit)
203
204
205 ggsave(' ../Plots/PCA plots/temp.png', p, height)
206
207 tiff(filename='C:/Users/Sandi/Dropbox/Projects/Liquid Gold/Manuscript/Manuscript A/
      Journal of Proteome Research/Figure S2.tif',
208      width=15*2/2.54, height=12*2/2.54, units='in', res=300,
209      pointsize=10)
210 grid.arrange(p.score, p, nrow=2)
211 grid.text('A', x=unit(0.7, 'cm'), y=unit(23.25, 'cm'), gp=gpar(cex=2))
212 grid.text('B', x=unit(0.7, 'cm'), y=unit(11.3, 'cm'), gp=gpar(cex=2))
213 dev.off()

```

Profile diversity analysis.R

```

1 #####
2 #Modified for Liquid Gold Manuscript A
3 #Profile diversity
4
5 #Setting working directory
6 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Analysis')
7
8 #loading required packages
9 require(RSQLite)
10 require(ggplot2)
11 require(plyr)
12 require(reshape2)
13 require(gridExtra)
14
15 #sourcing required functions

```

```

16 source(' ../Functions/function-data cleanup.R')
17 source(' ../Functions/function-descriptive statistics.R')
18 source(' ../Functions/function-scaling.R')
19 source(' ../Functions/function-extract legend.R')
20 source(' ../Functions/function-multiplot.R')
21 source(' ../Functions/function-heat map bin.R')
22 source(' ../Functions/function-defined binning.R')
23 source(' ../Functions/function-siegel tukey test.R')
24 source(' ../Functions/function-weighted compound class.R')
25
26
27 #setting working directory
28 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript A'
29 )
30 #Input data


---


31 # Choose driver
32 drv <- dbDriver('SQLite')
33
34 # Create connection (for SQLite, this is a file)
35 con <- dbConnect(drv, ' ../Data source files/data.db')
36
37 input <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
38                             q.hour, q.phase, q.dose, q.treatment, q.gen_treat, q.expt,
39                             d.variable, d.value, c.Type
40                             FROM Data d INNER JOIN Qualifier q
41                             ON d.file = q.file
42                             INNER JOIN Rank r
43                             ON d.variable = r.Compound AND q.expt = r.expt
44                             INNER JOIN Cmpd_type c
45                             ON r.Compound = c.Compound
46                             WHERE r.Confidence = 1 AND q.gen_treat = 'None'
47                             AND (q.run != 8 OR q.run IS NULL)
48                             AND d.variable NOT IN ('Xanthine', 'Glycerol')
49                             AND q.state = 'ss'
50                             AND q.donor != 'RP1'
51                             AND q.expt != 'RP2'
52                             ORDER BY q.sampleID;")
53
54 cmpd.info <- dbGetQuery(con, "SELECT c.Compound, c.Type, p.Pathway, m.main
55                                FROM Cmpd_type c INNER JOIN Pathway p
56                                ON c.Compound = p.Compound
57                                INNER JOIN Main_pathway m
58                                ON p.Pathway = m.path
59                                ORDER BY c.Compound")
60
61 dbDisconnect(con)
62
63 #data cleanup


---


64
65 #Performing data cleanup to remove zero value compounds
66 #also casts data into wide format
67 clean <- data_cleanup(input, convert.NA=TRUE)
68
69
70 #Ordering compound classes by abundance


---


71 rp <- input[input$donor == 'RP2',]
72 d <- input[input$donor != 'RP2',]
73
74 div.data <- clean
75 div.table <- as.data.frame(table(div.data$Type))
76 div.table <- div.table[order(div.table$Freq, decreasing=TRUE),]
77 var1.levels <- as.character(div.table$Var1)
78 div.data$Type <- factor(div.data$Type, levels=var1.levels, ordered=TRUE)
79
80 #Plotting profile diversity —> what times of compounds are the profiles made of

```

```

81 p <- ggplot(div.data, aes(x=donor, fill=Type)) +
82   geom_bar(stat='bin') +
83   theme_bw(14)
84 ggsave('./Plots/Profile diversity/profile diversity_order2.png', p, height=20,
85         width=15, units='cm')
86 #Abundance of compound scaled to concentrations
87 div <- dplyr::ddply(div.data, 'donor', custom_perc, 'Type')
88
89 p.div <- ggplot(div, aes(x=donor, y=class.avg, fill=Type)) +
90   geom_bar(stat='identity') +
91   theme_bw(14)
92
93 #export compount proportions
94 div.cast <- dcast(div, formula=Type~donor, value.var='perc')
95 write.csv(file='./Results/Profile diversity-weighted.csv', div.cast)
96
97 #Percentage of profiles is made of each compound class
98 -----
99  #(not scaled to concentration)
100 custom_table <- function(data){
101   table <- as.data.frame(table(data$Type))
102   prop <- table$Freq/sum(table$Freq)*100
103   type <- data.frame(donor=unique(data$donor), type=table$Var1, prop)
104   return(type)
105 }
106 table <- dplyr::ddply(clean, 'donor', custom_table)
107 type.data <- dcast(table, donor ~ type, value.var='prop')
108 write.csv(file='./Results/profile diversity2.csv', x=type.data)

```

Profile analysis.R

```

1 #####
2 #####
3 #Modified for Manuscript A
4 #Cross-experiment analysis
5
6 #Setting working directory
7 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Analysis')
8
9 #loading required packages
10 require(RSQLite)
11 require(ggplot2)
12 require(plyr)
13 require(reshape2)
14 require(gridExtra)
15 require(grid)
16 require(gridBase)
17
18 #sourcing required functions
19 source('../Functions/function-data cleanup.R')
20 source('../Functions/function-descriptive statistics.R')
21 source('../Functions/function-scaling.R')
22 source('../Functions/function-extract legend.R')
23 source('../Functions/function-multiplot.R')
24 source('../Functions/function-heat map bin.R')
25
26 #Input data
27 -----
28 # Choose driver
29 drv <- dbDriver('SQLite')
30
31 # Create connection (for SQLite, this is a file)
32 con <- dbConnect(drv, '../Data source files/data.db')
33
34 input <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
35                           q.hour, q.phase, q.dose, q.treatment, q.gen_treat,

```

```

35         q.expt, d.variable, d.value
36     FROM Data d INNER JOIN Qualifier q
37     ON d.file = q.file
38     INNER JOIN Rank r
39     ON d.variable = r.Compound AND q.expt = r.expt
40     WHERE r.Confidence = 1 AND q.gen_treat = 'None'
41     AND d.variable NOT IN ('Xanthine', 'Glycerol')
42     AND q.state = 'ss'
43     AND (q.run != 8 OR q.run IS NULL)
44     AND q.donor != 'RP1'
45     AND q.expt != 'RP2'
46     ORDER BY q.sampleID;" )
47
48 cmpd.info <- dbGetQuery(con, "SELECT c.Compound, c.Type, p.Pathway, m.main
49     FROM Cmpd_type c INNER JOIN Pathway p
50     ON c.Compound = p.Compound
51     INNER JOIN Main_pathway m
52     ON p.Pathway = m.path
53     ORDER BY c.Compound")
54
55 dbDisconnect(con)
56
57 #Setting working directory
58 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript A/
    Analysis')
59
60 #data cleanup
61
62 #Performing data cleanup to remove zero value compounds
63 #also casts data into wide format
64 clean <- data_cleanup(input, form.y=c('sampleID', 'donor', 'run', 'vessel',
65     'DPI', 'hour', 'phase', 'dose', 'treatment',
66     'gen_treat', 'expt'),
67     form.x='variable', convert.NA=TRUE)
68
69 #make sure all compounds have a value for each donor
70 data.cast <- dcast(clean, formula=sampleID+donor+run+vessel+DPI+hour+phase+dose+
71     treatment+gen_treat+expt~variable, value.var='value')
72
73 data.cast[is.na(data.cast)] <- 0
74
75 #melting back to long format
76 data.melt <- melt(data.cast, id.vars=c('sampleID', 'donor', 'run', 'vessel',
77     'DPI', 'hour', 'phase', 'dose', 'treatment',
78     'gen_treat', 'expt'))
79
80 #separating out qualifier and data
81 qualifier <- data.melt[,1:11]
82 x.data <- data.melt[,12:13]
83
84 #descriptive stats
85
86 #Getting descriptive stats for each donor
87 stat <- ddply(data.melt, 'donor', descriptive.stat, data.type='long', variable=
    variable,
88     qualifier=qualifier, scale.data=FALSE)
89
90 #density plot
91
92 #getting kernel density plot of values
93 p.dens <- ggplot(stat[stat$avg >= 0.1 & stat$avg <=20,]) +
94     geom_density(aes(x=avg)) +
95     theme_bw(16)
96
97 plot(density(stat$avg))

```

```

98 plot(density(stat$avg))
99 png(file='../Plots/bin_density.png', height=10, width=15, units='cm', res=240)
100 plot(density(stat$avg), main='', xlab='Mean Concentration (mM)',
101       sub='bw=0.03089')
102 dev.off()
103
104 temp <- stat$avg[stat$avg >=0.1 & stat$avg <=20]
105 plot(density(temp)) #bw=0.135
106
107 png(file='../Plots/bin_density-inset.png', height=7.5, width=10, units='cm', res
      =240)
108 plot(density(stat$avg[stat$avg >=0.1 & stat$avg <=20]),
109       main='', xlab='', ylab='', cex=10, xaxt='n')
110 axis(1, at=seq(0,20,5), labels=c(0.1, seq(5,20,5)))
111 dev.off()
112
113 vp <- viewport(x=0.55, y=0.3, width=0.95, height=0.95)
114 vp.inset <- viewport(x=0.55, y=0.8, width=unit(12, 'cm'),
115                      height=unit(8.5, 'cm'))
116 tiff(filename='C:/Users/Sandi/Dropbox/Projects/Liquid Gold/Manuscript/Manuscript A/
      Journal of Proteome Research/Figure S1.tif',
117       width=7, height=5.5, units='in', res=300,
118       pointsize=10)
119 plot.new()
120 pushViewport(vp)
121 plot(density(stat$avg), main='', xlab='Mean Concentration (mM)',
122       sub='bw=0.03089')
123 pushViewport(vp.inset)
124 par(new=TRUE, fig=gridFIG())
125 plot(density(stat$avg[stat$avg >=0.1 & stat$avg <=20]),
126       main='', xlab='', ylab='', cex=10, xaxt='n')
127 axis(1, at=seq(0,20,5), labels=c(0.1, seq(5,20,5)))
128 upViewport()
129 grid.lines(x=c(0.85,0.26), y=c(0.64,0.375))
130 grid.lines(x=c(0.298,0.062), y=c(0.64,0.45))
131 dev.off()
132
133 #Binning


---


134 #Assigning data to bins
135 bin.data <- bin.heatmap(stat, bin=c(20,10,1,0.1,0), 'variable',
136                            order.data=stat, order.var='variable', order.value='avg')
137
138 uni.list <- bin.data[['uni.list']]
139 ordered.data <- bin.data[['ordered.data']]
140 ordered.data$donor <- factor(ordered.data$donor, levels=c('RP2','A','B','C'))
141
142 #Padded variable names
143 ordered.data$pad <- sprintf('%022s', as.character(ordered.data$variable))
144
145 pad.list <- list()
146 for(i in 1:length(uni.list)) {
147   pad.temp <- sprintf('%022s', as.character(uni.list[[i]]))
148   pad.list[[i]] <- pad.temp
149 }
150
151 pad.order <- c()
152 for(i in 1:length(pad.list)) {
153   temp <- as.character(pad.list[[i]])
154   pad.order <- c(pad.order, temp)
155 }
156
157 #plotting heat maps


---


158 #plotting heat maps by bin
159 plot.list <- list()
160 legend.list <- list()

```

```

162 fig.list <- list()
163 fig.leg <- list()
164
165
166 fig.width <- c(7.7, 7.5, 7.8, 8.8, 9.4)
167 fig.height <- c(2.65, 3.1, 4.7, 13.5, 10.9)
168 margin <- list(c(0, 0.2, 0.5, 0.8),
169               c(0, 0.2, 0.5, 0.8),
170               c(0, 0.2, 0.3, 0.4),
171               c(0, 0.2, 0, 0),
172               c(0, 0.2, 0, 0))
173
174 for(i in 1:length(uni.list)) {
175   plot.data <- ordered.data[ordered.data$variable %in% uni.list[[i]],]
176
177   #setting compound order
178   plot.data$pad <- factor(as.character(plot.data$pad),
179                           levels=rev(pad.order[pad.order %in% pad.list[[i]]]),
180                           ordered=TRUE)
181
182   plot.list[[i]] <- ggplot(plot.data, aes(x=donor, y=pad)) +
183     geom_tile(aes(fill=avg), colour='grey90') +
184     scale_fill_gradient(low='grey90', high='grey0',
185                        name='Mean\nConcentration (mM)') +
186     scale_x_discrete(labels=c('MET-2', 'A', 'B', 'C')) +
187     theme_bw(10) +
188     xlab('Donor') +
189     theme(legend.position=c(0.52, 0.5), legend.direction='horizontal',
190           legend.key.width=unit(1, 'cm'),
191           legend.key.height=unit(0.3, 'cm'),
192           legend.text=element_text(size=7.5),
193           legend.title=element_text(size=7.5),
194           axis.title.y=element_blank(),
195           plot.margin=unit(margin[[i]], 'cm'))
196
197   legend.list[[i]] <- g_legend(plot.list[[i]])
198
199   fig.list[[i]] <- grid.arrange(legend.list[[i]],
200                                 plot.list[[i]] + theme(legend.position='none'),
201                                 heights=c(0.1, 0.9))
202
203   #save as panels
204   if(FALSE){
205     png(filename=paste('.. /Plots/Profiles/profile4_bin', i, '.png', sep=''),
206         height=fig.height[i], width=fig.width[i], units='cm', res=240, pointsize
207         =10)
208     fig.list[[i]] <- grid.arrange(plot.list[[i]] + theme(legend.position='none'))
209     dev.off()
210
211     png(filename=paste('.. /Plots/Profiles/profile4_bin_leg', i, '.png', sep=''),
212         height=1, width=leg.width[i], units='cm', res=240, pointsize=10)
213     fig.leg[[i]] <- grid.arrange(legend.list[[i]])
214     dev.off()
215   }
216 }
217
218 #Setting plot dimensions
219
220 w <- list(c(10, 7.7), #1
221          c(10, 7.5), #2
222          c(10, 7.8), #3
223          c(10, 8.8), #4
224          c(10, 9.4)) #5
225
226 h <- list(c(1, 2.5), #1

```

```

227         c(1,2.9),      #2
228         c(1,4.5),      #3
229         c(1,13.5),     #4
230         c(1,10.3))     #5
231
232 #EPS
233 tiff(filename='C:/Users/Sandi/Dropbox/Projects/Liquid Gold/Manuscript/Manuscript A/
      Journal of Proteome Research/Figure 3.tif',
234       width=7, height=sum(c(h[[4]],h[[3]])/2.54), units='in', res=300,
235       pointsize=11)
236 grid.arrange(arrangeGrob(legend.list[[1]],
237                          plot.list[[1]] + theme(legend.position='none'),
238                          legend.list[[2]],
239                          plot.list[[2]] + theme(legend.position='none'),
240                          legend.list[[5]],
241                          plot.list[[5]] + theme(legend.position='none'),
242                          ncol=1, heights=c(h[[1]],h[[2]],h[[5]]),
243                          widths=c(w[[1]],w[[2]],w[[5]])),
244       arrangeGrob(legend.list[[3]],
245                  plot.list[[3]] + theme(legend.position='none'),
246                  legend.list[[4]],
247                  plot.list[[4]] + theme(legend.position='none'),
248                  ncol=1, heights=c(h[[3]],h[[4]]),
249                  widths=c(w[[3]],w[[4]])),
250       ncol=2)
251 grid.text('A', x=unit(0.25,'cm'), y=unit(19.45, 'cm'),gp=gpar(cex=1.5))
252 grid.text('B', x=unit(0.25,'cm'), y=unit(15.7, 'cm'),gp=gpar(cex=1.5))
253 grid.text('E', x=unit(0.25,'cm'), y=unit(11.52, 'cm'),gp=gpar(cex=1.5))
254 grid.text('C', x=unit(9.1,'cm'), y=unit(19.45, 'cm'),gp=gpar(cex=1.5))
255 grid.text('D', x=unit(9.1,'cm'), y=unit(14, 'cm'),gp=gpar(cex=1.5))
256 dev.off()
257
258 #Mann Whitney test
259
260 #seeing what compounds differ b/w donor and RPs
261 mw_custom <- function(data, value='value') {
262   d <- data[data$donor %in% c('A','B','C'),]
263   rp <- data[data$donor == 'RP2',]
264   out <- c()
265   test.result=wilcox.test(d[,value], rp[,value])
266
267   output <- rp
268
269   output$p.val <- test.result$p.value
270   output$stat <- test.result$statistic
271
272   out <- rbind(out, output)
273   return(out)
274 }
275
276
277 mw.result <- ddply(stat, 'variable', mw_custom)
278
279 sig <- mw.result[mw.result$p.val <= 0.05,]
280
281
282 sig.bin <- uni.list
283
284 #Putting significant compounds in bin order
285 sig.order <- list()
286 for(i in 1:length(sig.bin)) {
287   sig.data <- sig[sig$variable %in% sig.bin[[i]],]
288   sig.data <- unique(sig.data[,c('variable','p.val')])
289
290   type <- cmpd.info[cmpd.info$Compound %in% unique(sig.data$variable),]
291   for(j in unique(type$Compound)){
292     type$p.val[type$Compound == j] <- sig.data$p.val[sig.data$variable == j]

```

```

293   #print('end1 ')
294 }
295
296 #type$Compound <- factor(type$Compound,
297 #                           levels=rev(levels(sig.bin[[i]])), order=TRUE)
298 #type <- type[order(type$Compound),]
299 type <- type[order(type$P.val, decreasing=FALSE),]
300
301 sig.order[[i]] <-type
302 #print('end2 ')
303
304 }
305
306 #Putting all in one dataframe — Compound Type
307 sig.all.cmpd <- do.call(rbind.data.frame, sig.order)
308 sig.all.cmpd <- unique(sig.all.cmpd[,c(1,5,2)])
309 table(sig.all.cmpd$Type)
310
311 #Putting all in one dataframe — Pathway
312 sig.all.path <- do.call(rbind.data.frame, sig.order)
313 sig.all.path <- unique(sig.all.path[,c(1,5,3)])
314 table(sig.all.path$Pathway)
315
316 #write.csv(file = '../Results/donor-pvalue2.csv', sig.all)
317
318 #exporting each bin separately
319 for(i in 1:length(sig.order)) {
320   sig.write <- sig.order[[i]]
321   #write.csv(file=paste('../Results/donor-pvalue',i,'_2.csv',sep=''), sig.write)
322 }
323
324 #Seeing how these compounds differ across donors—————
325 d <- sig.all
326
327 out <- data.frame(A=NULL, B=NULL, C=NULL, RP2=NULL)
328 for(i in d$Compound){
329   avg <- unique(stat[stat$variable == i, c('donor','avg')])
330   tmp <- as.data.frame(t(avg))
331   val <- tmp[2,]
332   colnames(val) <- as.character(as.matrix(tmp[1,]))
333   val$variable <- i
334   out <- rbind(out, val)
335 }
336
337 d <- cbind(d, out)
338
339 #Getting mean concentrations of each group—————
340 a <- stat[stat$variable %in% sig.all$Compound,]
341
342 rp <- a[a$donor == 'RP2',]
343 dn <- a[a$donor %in% c('A','B','C'),]
344
345 rp$donor.type <- 'rp'
346 dn$donor.type <- 'd'
347
348 rp <- ddply(rp, 'variable', mutate, rp.avg=mean(unique(avg)))
349 dn <- ddply(dn, 'variable', mutate, dn.avg=mean(unique(avg)))
350
351 rp2 <- unique(rp[,c('donor.type','variable','rp.avg')])
352 dn2 <- unique(dn[,c('donor.type','variable','dn.avg')])
353
354 b <- cbind(rp2[c('variable','rp.avg')], 'dn.avg'=dn2$dn.avg)
355
356 source('../Functions/defined binning.R')
357 bin.info <- defined_binning()
358 b$bin <- NULL
359 for(i in b$variable)b$bin[b$variable == i] <- bin.info$bin[bin.info$variable==i]
360

```



```

361 h.rp <- b[b$rp.avg > b$dn.avg,]
362 l.rp <- b[b$rp.avg < b$dn.avg,]
363
364 #Visualize in plot form
365 d.melt <- melt(d, id.vars=c('Compound','variable','Type','p.val'),
366               measure.vars=c('A','B','C','RP2'), variable.name='donor')
367 d.melt$value <- as.numeric(d.melt$value)
368
369 plot.var <- unique(d.melt$variable)
370 plot.var <- factor(plot.var, levels=unlist(uni.list))
371 plot.var <- plot.var[order(plot.var)]
372 plot.var <- as.character(plot.var)
373
374 plot.function <- function(.x, dat, var, variable) {
375   plot.data <- dat[dat[,variable]==var[.x],]
376   p <- ggplot(plot.data, aes(x=donor, y=variable)) +
377     geom_tile(aes(fill=value)) +
378     geom_text(aes(fill=value, label=round(plot.data$value,3))) +
379     scale_fill_gradient(low='white', high='blue') +
380     theme_bw(14) +
381     theme(axis.title.y=element_blank())
382   return(p)
383 }
384 p <- lapply(1:length(plot.var), plot.function,
385            dat=d.melt, var=plot.var, variable='variable')
386
387 png('../Plots/test.png', height=100, width=15, units='cm', res=240)
388 do.call(grid.arrange, c(p, ncol=1))
389 dev.off()
390 #pathway analysis
391 path.table <- as.data.frame(table(data$Pathway))
392 path.table <- path.table[order(path.table$Freq, decreasing=TRUE),]
393 #write.csv(file='../Results/donor-pathway.csv', path.table)
394
395 main.table <- as.data.frame(table(data$main))
396 main.table <- main.table[order(main.table$Freq, decreasing=TRUE),]
397 #write.csv(file='../Results/donor-main pathway.csv', main.table)
398
399 data$main.count <- NA
400 for(i in unique(main.table$Var1)) {
401   data$main.count[data$main == i] <- main.table$Freq[main.table$Var1 == i]
402 }
403
404 data$path.count <- NA
405 for(i in unique(path.table$Var1)){
406   data$path.count[data$Pathway == i] <- path.table$Freq[path.table$Var1 == i]
407 }
408
409 data <- data[order(-data$path.count),]
410 path.level <- unique(data$Pathway)
411 data$Pathway <- factor(data$Pathway, levels=rev(path.level), ordered=TRUE)
412 cmpd.level <- unique(data$Compound)
413 data$Compound <- factor(data$Compound, levels=cmpd.level, ordered=TRUE)
414
415 p.path <- ggplot(data[data$main %in% c('Amino acids', 'Organic acids',
416                                     'Vitamins', 'Central Metabolic Pathways',
417                                     'Lipids', 'Purine', 'Pyrimidine', 'Fatty acids')
418                  ],
419                aes(x=main, fill=Pathway)) +
420   geom_bar(stat='bin', colour='black') +
421   theme_bw(14) +
422   theme(axis.text.x=element_text(hjust=0, vjust=1, angle=-45))

```

Biomarker analysis.R

```

1 #####
2 #Analysis for biomarkers unique to each donor
3 #Created for Liquid Gold Manuscript A

```

```

4  #(not in regular Liquid Gold manuscript R analyses)
5
6  #Setting working directory
7 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Analysis')
8
9  #loading required packages
10 require(RSQLite)
11 require(ggplot2)
12 require(plyr)
13 require(reshape2)
14 require(gridExtra)
15 require(stringr)
16
17  #sourcing required functions
18 source(' ../Functions/function-data cleanup.R')
19 source(' ../Functions/function-descriptive statistics.R')
20 source(' ../Functions/function-scaling.R')
21 source(' ../Functions/function-extract legend.R')
22 source(' ../Functions/function-multiplot.R')
23 source(' ../Functions/function-heat map bin.R')
24 source(' ../Functions/function-defined binning.R')
25 source(' ../Functions/function-siegel tukey test.R')
26
27
28  #setting working directory
29 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript A'
30 )
31  #Input data


---


32  # Choose driver
33 drv <- dbDriver('SQLite')
34
35  # Create connection (for SQLite, this is a file)
36 con <- dbConnect(drv, '../Data source files/data.db')
37
38 input <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
39 q.hour, q.phase, q.dose, q.treatment, q.gen_treat, q.expt,
40 d.variable, d.value, c.Type
41 FROM Data d INNER JOIN Qualifier q
42 ON d.file = q.file
43 INNER JOIN Rank r
44 ON d.variable = r.Compound AND q.expt = r.expt
45 INNER JOIN Cmpd_type c
46 ON r.Compound = c.Compound
47 WHERE r.Confidence = 1 AND q.gen_treat = 'None'
48 AND (q.run != 8 OR q.run IS NULL)
49 AND d.variable NOT IN ('Xanthine', 'Glycerol')
50 AND q.state = 'ss'
51 AND q.donor != 'RP1'
52 AND q.expt != 'RP2'
53 ORDER BY q.sampleID;")
54
55 cmpd.info <- dbGetQuery(con, "SELECT c.Compound, c.Type, p.Pathway, m.main
56 FROM Cmpd_type c INNER JOIN Pathway p
57 ON c.Compound = p.Compound
58 INNER JOIN Main_pathway m
59 ON p.Pathway = m.path
60 ORDER BY c.Compound")
61
62 dbDisconnect(con)
63
64  #data cleanup


---


65
66  #Performing data cleanup to remove zero value compounds
67  #also casts data into wide format
68 clean <- data_cleanup(input, convert.NA=TRUE)

```

```

69
70 #make sure all compounds have a value for each donor
71 data.cast <- dcast(clean, formula=sampleID+donor+run+vessel+DPI+hour+phase+dose+
72   treatment+gen_treat+expt~variable, value.var='value')
73
74 data.cast[is.na(data.cast)] <- 0
75
76 #melting back to long format
77 data.melt <- melt(data.cast, id.vars=c('sampleID','donor','run','vessel',
78   'DPI','hour','phase','dose','treatment',
79   'gen_treat','expt'))
80
81 #descriptive stats
82
83 #Getting descriptive stats for each donor
84 stat <- ddply(data.melt, 'donor', descriptive.stat, data.type='long', variable='
85   variable',
86   qualifier=NULL, scale.data=FALSE)
87
88 #Metabolite markers
89 #markers identified visually
90 coi <- c('3-Phenylpropionate','Leucine','Valine','Isoleucine',
91   'Methylamine','Methionine','Phenylalanine','Methanol','Succinate',
92   'Dimethylamine', ##Donor A
93
94   'Pimelate','Sebacate','4-Hydroxyphenylacetate','Acetate','Ethanol',
95   'Valerate','Lactate', ##Donor B
96
97   'beta-Alanine','Pyruvate','Glucose', ##Donor C, low in Lactate too
98
99   'Tartrate','Urea','Glutamate','Pyroglutamate','Asparagine','Glycolate',
100   'Choline','Thymine','Formate','Isovalerate','Isobutyrate','p-Cresol',
101   'Lysine') ##RP2
102 coi <- factor(coi, levels=coi, ordered=TRUE)
103
104 #Seeing if same compounds picked by anova
105
106 #tukey fails because too few data points
107 source('../Functions/function-tukey.R')
108 source('../Functions/function-anova-custom.R')
109
110 anova <- ddply(stat, 'variable', aov-custom, group='donor')
111 s.anova <- anova[anova$p.val <= 0.05 & !is.na(anova$p.val),]
112
113 tukey <- ddply(stat, 'variable', tukey-custom, group='donor')
114 s.tukey <- tukey[tukey$p.val <= tukey$bon.val | is.na(tukey$p.val),]
115
116 #write.csv(file='Results/tukey-sig.csv', s.tukey)
117
118 #see if bonferroni correction picked anything not in coi
119 bon.tukey <- s.tukey[!s.tukey$variable %in% coi,]
120
121 #Check if bonferroni correction cut out any coi
122 cut.tukey <- coi[!coi %in% unique(s.tukey$variable)]
123
124 cut.tukey <- tukey[tukey$variable %in% cut.tukey,]
125
126 ##bonferroni cut out compounds found only in one donor AND had large spread
127 ##these compounds still need to be discussed
128
129 write.csv(file='../Results/profile-diversity-tukey-bonferroni.csv', s.tukey)
130
131 #getting compounds where one group is significantly different from every single
132   group
133 marker-custom <- function(data) {
134   donor <- c(as.character(data$g1), as.character(data$g2))
135   t <- as.data.frame(table(donor))
136   t$variable <- unique(data$variable)
137   marker <- t[t$Freq == 3,]

```

```

132 out <- marker[,c('variable','donor')]
133 return(out)
134 }
135
136 marker.tukey <- ddply(s.tukey, 'variable', marker_custom)
137
138 #Getting markers that only have one donor that is different from every single other
group
139 table.marker <- as.data.frame(table(as.character(marker.tukey$variable)))
140 uni.marker <- as.character(table.marker$Var1[table.marker$Freq == 1])
141
142 marker.tukey2 <- marker.tukey[marker.tukey$variable %in% uni.marker,]
143
144 #compare to those done by visual inspection
145
146
147 write.csv(file='./Results/profile_diversity-marker_tukey.csv', marker.tukey)
148 write.csv(file='./Results/profile_diversity-marker2_tukey.csv', marker.tukey2)
149 #Seeing if same compounds picked by t-test

```

```

150 #t-test with Bonferroni correction
151 t_custom <- function(data, group) {
152
153   grps <- unique(data[,group])
154   pair <- combn(grps,2)
155
156   out <- c()
157   for(i in 1:ncol(pair)) {
158     grp1 <- data[data[,group] == pair[1,i],]
159     grp2 <- data[data[,group] == pair[2,i],]
160     grp <- rbind(grp1, grp2) #both groups being compared
161
162     result <- t.test(grp1$value, grp2$value)
163     p.val <- result$p.value
164     statistic <- result$statistic
165     output <- data.frame(g1=pair[1,i], g2=pair[2,i], stat=statistic, p.val=p.val)
166
167     out <- rbind(out, output)
168   }
169   return(out)
170 }
171
172 t.result <- ddply(stat,'variable', t_custom, group='donor')
173
174 pair <- combn(unique(stat$donor),2)
175 bonferroni <- 0.05/ncol(pair)
176
177 sig.t <- t.result[t.result$p.val <= bonferroni & !is.na(t.result$p.val),]
178
179 #see if bonferroni correction picked anything not in coi
180 bon <- sig.t[!sig.t$variable %in% coi,]
181
182 #Check if bonferroni correction cut out any coi
183 cut <- coi[!coi %in% unique(sig.t$variable)]
184 cut <- t.result[t.result$variable %in% cut,]
185
186 other0s <- t.result[t.result$variable %in% c('Glucose'),]
187 ##bonferroni cut out compounds found only in one donor AND had large spread
188 ##these compounds still need to be discussed
189
190 write.csv(file='./Results/profile_diversity-t_bonferroni.csv', sig.t)
191
192 #getting compounds where one group is significantly different from ever single
group
193 marker_custom <- function(data) {
194   donor <- c(as.character(data$g1), as.character(data$g2))
195   t <- as.data.frame(table(donor))
196   t$variable <- unique(data$variable)

```

```

197   print(t)
198   marker <- t[t$Freq == 3,]
199   out <- marker[,c('variable', 'donor')]
200   return(out)
201 }
202
203 marker <- ddply(sig.t, 'variable', marker_custom)
204
205 #Getting markers that only have one donor that is different from every single other
group
206 table.marker <- as.data.frame(table(as.character(marker$variable)))
207 uni.marker <- as.character(table.marker$Var1[table.marker$Freq == 1])
208
209 marker2 <- marker[marker$variable %in% uni.marker,]
210
211 #compare to those done by visual inspection
212 cut.sig <- coi[!coi %in% unique(sig.t$variable)] #cut by bonferroni
213 t.cutsig <- t.result[t.result$variable %in% cut.sig,]
214 cut.m1 <- coi[!coi %in% unique(marker$variable)] #cut by marker definition
215 t.cut1 <- sig.t[sig.t$variable %in% cut.m1,]
216 cut.m2 <- coi[!coi %in% unique(marker2$variable)] #cut by marker for one donor
217 t.cut2 <- sig.t[sig.t$variable %in% cut.m2,]
218
219 write.csv(file='./Results/profile_diversity-marker.csv', marker)
220 write.csv(file='./Results/profile_diversity-marker2.csv', marker2)
221 #Getting plot data and order-----
222
223 p2.data <- stat[stat$variable %in% coi,]
224 p2.data$variable <- factor(p2.data$variable, levels=levels(coi))
225
226 pattern <- '(*Hydroxy|henyl)(?!alanine)(.*)'
227 match <- '\\1-\\n\\2'
228
229 new.values <- as.character(gsub(pattern, match, p2.data$variable, perl=TRUE))
230 new.levels <- as.character(gsub(pattern, match, levels(p2.data$variable), perl=TRUE
  ))
231
232 p2.data$wrap <- factor(new.values, levels=new.levels)
233
234 #Plotting out data values-----
235 p2 <- ggplot(p2.data, aes(x=donor, y=value)) +
236   geom_errorbar(aes(ymax=ymax, ymin=ymin)) +
237   geom_point(size=1.5) +
238   facet_wrap(~wrap, scales='free', ncol=5) +
239   theme_bw(10) +
240   scale_x_discrete(labels=c('A', 'B', 'C', 'MET-2')) +
241   ylab('Concentration (mM)') +
242   xlab('Donor')
243 ggsave('./Plots/Profile_diversity/markers3.png', p2,
244   height=24, width=7*2.54, units='cm')
245
246 tiff(file='C:/Users/Sandi/Dropbox/Projects/Liquid_Gold/Manuscript/Manuscript A/
  Journal of Proteome Research/Figure 4.tif',
247   height=24/2.54, width=7, units='in', res=300,)
248 print(p2)
249 dev.off()

```

C.5 Analysis of Perturbations of Human Fecal Microbiota Propagated In Vitro

PCA analysis.R

```

1 #####
2 #Modified for Manuscript B

```

```

3
4 #Analysis to compare samples across projects
5
6 #loading requiried packages
7 require(RSQLite)
8 require(ggplot2)
9
10 #sourcing required functions
11 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions/')
12 source('function-pca.R')
13 source('function-plotting 2D score.R')
14 source('function-plotting 2D loading.R')
15 source('function-pca diagnostics.R')
16 source('function-data cleanup.R')
17 source('function-scaling.R')
18
19 #setting working directory
20 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript B/
21 )
22 #
23 #Input data
24 # Choose driver
25 drv <- dbDriver('SQLite')
26
27 # Create connection (for SQLite, this is a file)
28 con <- dbConnect(drv, '../Data source files/data.db')
29
30 input <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
31                          q.hour, q.phase, q.dose, q.treatment, q.gen_treat, q.expt,
32                          d.variable, d.value
33                          FROM Data d INNER JOIN Qualifier q
34                          ON d.file = q.file
35                          INNER JOIN Rank r
36                          ON d.variable = r.Compound AND q.expt = r.expt
37                          WHERE r.Confidence = 1 AND q.gen_treat = 'None'
38                          AND d.variable NOT IN ('Xanthine', 'Glycerol')
39                          AND q.state = 'ss'
40                          AND q.donor NOT IN ('C', 'RP1')
41                          AND q.expt != 'RP2'
42                          ORDER BY q.sampleID;")
43
44 input2 <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
45                          q.hour, q.phase, q.dose, q.treatment, q.gen_treat, q.expt,
46                          d.variable, d.value
47                          FROM Data d INNER JOIN Qualifier q
48                          ON d.file = q.file
49                          INNER JOIN Rank r
50                          ON d.variable = r.Compound AND q.expt = r.expt
51                          WHERE r.Confidence = 1
52                          AND d.variable NOT IN ('Xanthine', 'Glycerol')
53                          AND q.state = 'ss'
54                          AND q.donor NOT IN ('C', 'RP1')
55                          AND q.expt != 'RP2'
56                          ORDER BY q.sampleID;")
57 dbDisconnect(con)
58
59 #setting working directory
60 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript B/
61 Analysis')
62 #
63 #Performing cleanup function to remove DSS, contaminants and compounds of
64 #zero concentration
65
66 clean.data <- data_cleanup(input2, 'sampleID', 'variable', 'value',
67                             form.y=c('sampleID', 'donor', 'run', 'vessel', 'DPI',
68                                     'hour', 'phase', 'dose', 'treatment', 'gen_treat',

```

```

69         'expt'),
70         form.x='variable', contaminants=NULL, export=FALSE)
71
72 #-----
73 #converting factors to appropriate classes
74 clean.data$run <- as.character(clean.data$run)
75 clean.data$gen_treat <- factor(clean.data$gen_treat,
76                               levels=c('None', 'Norepinephrine', 'RePOOPulate',
77                                         'Clindamycin+RePOOPulate', 'Clindamycin',
78                                         'Vancomycin'))
79 #-----
80 #performing the PCA
81 pca.data = perform_pca(clean.data, sampleID='sampleID',
82                        variable='variable', value='value',
83                        form.y=c('sampleID', 'donor', 'run', 'vessel', 'DPI',
84                                'hour', 'phase', 'dose', 'treatment', 'gen_treat',
85                                'expt'),
86                        form.x='variable',
87                        scale='UV', center=TRUE, convert.NA=TRUE)
88 #-----
89 #performing diagnostics on PCA
90 diagnostic.results = scree_plot(pca.data)
91
92 #unpacking the diagnostic results
93 diagnostic.number = diagnostic.results[['number']]
94 diagnostic.explained = diagnostic.results[['explained']]
95 diagnostic.plot = diagnostic.results[['plot']]
96 diagnostic.plot = diagnostic.plot + ggtitle('Scree Plot (UVN)')
97
98 #printing out diagnostic results
99 print(diagnostic.number)
100 print(diagnostic.explained)
101 print(diagnostic.plot)
102 #-----
103 #plotting PCA-processed data
104 #Score plot 2D
105 # Note: Arguments for score plot function are
106 #       (pca_data, qualifier_data, PCs=1:4, colour=NULL, shape=NULL, size=NULL,
107 #       label=NULL, title=NULL, legend.ori='horizontal', black-white=FALSE)
108
109 score_plot = score_plot_2d(pca.data, PCs=1:4,
110                           colour='gen_treat', shape='donor',
111                           legend.ori='horizontal', outline=FALSE,
112                           black-white=FALSE)
113
114 p.score <- score_plot[['p']]
115 collected.data <- score_plot[['collected.data']]
116
117 ellipse_custom <- function(d) {
118
119   #Calculating ellipse for 95% confidence
120   ellipseSD = cor(d$x.value, d$y.value)
121   ellipseDf = as.data.frame(ellipse(ellipseSD, scale=c(sd(d$x.value),
122                                                         sd(d$y.value)),
123                                     centre=colMeans(d[,c('x.value', 'y.value')]),
124                                     t=2))
125   return(ellipseDf)
126 }
127
128 e.path <- ddply(collected.data, 'view', ellipse_custom)
129
130
131 #Drawing an ellipse for 95% confidence
132 p.score = p.score +
133   # geom_path(data=e.path, aes(x=x, y=y), linetype=2, alpha=0.5) +
134   # scale_colour_manual(values=c('black', 'darkorange1', 'cornflowerblue',
135   #                               'black', 'chartreuse', 'blue', 'grey30'),
136   #
137   guide_legend(title='Run'),

```

```

137 # labels=c('NA', '1', '2', '3', '4', '5', '6')) +
138 scale_shape_manual(name='Donor', labels=c('A', 'B', 'MET-2'), values=c(16,7,5)) +
139 scale_colour_manual(values=c('cyan4', 'darkorchid4', 'chartreuse',
140 'darkorange1', 'deeppink1', 'blue', 'grey30'),
141 labels=c('None', 'Norepinephrine', 'MET-1',
142 'Clindamycin+MET-1', 'Clindamycin',
143 'Vancomycin'),
144 guide_legend(title='Treatment',
145 )) +
146 guides(col=guide_legend(byrow=TRUE, nrow=2))
147
148 ggsave('.. /Plots/PCA plots/treatment-score.png', p.score, height=15, width=25,
149 units='cm')
150 #
151
152 #Loading plot 2D
153 # Note: Arguments for loading plot function are (pca_data, PCs=1:4, label=NULL,
154 title=NULL,
155 # view_mode='compressed')
156 #If you want all compound labels shown, enter all_paths in label argument
157 all_paths = rownames(pca.data$rotation)
158 #If you want all and only compound within the subset of observations processed in
159 the PCA function, can create a temporary group called 'a' to enter in for the
160 label argument
161 #a = pathway_list$Carb[pathway_list$Organic %in% rownames(pca_processed_data$
162 rotation)]
163
164 loading_plot = loading_plot_2d(pca.data, PCs=1:4, label=all_paths, outline=TRUE,
165 title=NULL, view_mode = 'compressed')
166
167 p.load <- loading_plot[['p']]
168 load.data <- loading_plot[['collected.data']]
169
170 ggsave('.. /Plots/PCA plots/treatment-loading.png', p.load,
171 height=14.5*1.5, width=25*2, units='cm')

```

Treatment analysis.R

```

1 #####
2 //
3 #Modified for Manuscript B
4 #Antibiotic treatment analysis
5 #Heat map
6
7 #loading required packages
8 require(ggplot2)
9 require(RODBC)
10 require(RSQLite)
11 require(gridExtra)
12 require(scales)
13
14 #sourcing required functions
15 setwd('c:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions/')
16 source('function-descriptive statistics.R')
17 source('function-data cleanup.R')
18 source('function-extract legend.R')
19 source('function-bin function.R')
20
21 #setting working directory
22 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript B/
23 Analysis')
24 source('.. /Functions/function-defined binning.R')
25 #

```



```

25 #Input data
26 con <- odbcConnect('MySQLDSN')
27
28 input <- sqlQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
29       q.hour, q.phase, q.dose, q.treatment, q.gentreat, q.expt,
30       d.variable, d.value
31       FROM liquidgold.data d INNER JOIN liquidgold.qualifier q
32       ON d.file = q.file
33       INNER JOIN liquidgold.rank r
34       ON d.variable = r.Compound AND q.expt = r.expt
35       WHERE r.Confidence = 1
36       AND (q.run != 8 OR q.run IS NULL)
37       AND d.variable NOT IN ('Xanthine', 'Glycerol')
38       AND q.state = 'ss'
39       AND q.donor NOT IN ('C', 'RP1')
40       AND q.expt != 'RP2'
41       ORDER BY q.sampleID;")
42
43 odbcClose(con)
44
45 #
46 #Performing data cleanup to remove zero value compounds
47 #also casts data into wide format
48 clean <- data_cleanup(input, form.y=c('sampleID', 'donor', 'run', 'vessel',
49       'DPI', 'hour', 'phase', 'dose', 'treatment',
50       'gentreat', 'expt'),
51       form.x='variable', convert.NA=TRUE)
52
53 #make sure all compounds have a value for each donor
54 data.cast <- dcast(clean, formula=sampleID+donor+run+vessel+DPI+hour+phase+dose+
55       treatment+gentreat+expt~variable, value.var='value')
56
57 #Assigning 0 to where compound was not profiled
58 data.cast[,12:ncol(data.cast)][is.na(data.cast[,12:ncol(data.cast)])] <- 0
59
60 #melting back to long format
61 data.melt <- melt(data.cast, id.vars=c('sampleID', 'donor', 'run', 'vessel',
62       'DPI', 'hour', 'phase', 'dose', 'treatment',
63       'gentreat', 'expt'))
64 #
65 #Getting descriptive stats for compounds within gentreat for each donor
66 stat <- ddply(data.melt, c('donor', 'gentreat'), descriptive.stat, data.type='long',
67       variable='variable')
68
69 #t test
70
71 #testing if the change from non-treated to treated is significant usint t-test
72 t_custom <- function(data, gentreat) {
73
74   none <- data[data$gentreat == 'None',]
75   out <- none
76   out$p.val <- NA
77
78   for(i in gentreat){
79
80     treat <- data[data$gentreat == i,]
81
82     #when treated samples have more than 1 value, can use comparison of 2 means
83     if(length(unique(treat$sampleID)) > 1) {
84       p.value=t.test(none$value, treat$value)$p.value
85       output <- treat
86       output$p.val <- p.value
87
88       out <- rbind(out, output)
89     }
90
91     #when treated samples have only 1 value, do 1-sample t test, where treated

```

```

92     value
93     #is the true mean
94     else if (length(unique(treat$sampleID)) == 1) {
95         x.bar <- mean(none$value)
96         y.value <- treat$value
97         x.sd <- sd(none$value)
98         n <- length(unique(none$sampleID))
99
100         t.stat <- (x.bar - y.value)/(x.sd/sqrt(n))
101         p.value <- 2*pt(-abs(t.stat), n-1)
102
103         output <- treat
104         output$p.val <- p.value
105
106         out <- rbind(out, output)
107     }
108     #when the gentreat is not in that donor
109     else if (length(unique(treat$sampleID)) < 1) {
110         out <- out
111     }
112 }
113 return(out)
114 }
115
116 t.result <- ddply(stat, c('donor', 'variable'),
117                   t_custom, c('Clindamycin', 'RePOOPulate',
118                               'Clindamycin+RePOOPulate', 'Norepinephrine', '
119                               Vancomycin'))
120
121 #write.csv(file='check_tresult.csv', t.result)
122
123 #Mann Whitney test
124
125 #testing if the change from non-treated to treated is significant using mann
126 whitney test
127
128 mw_custom <- function(data, gentreat) {
129
130     none <- data[data$gentreat == 'None',]
131     out <- none
132     out$p.val <- NA
133
134     for(i in gentreat){
135
136         treat <- data[data$gentreat == i,]
137
138         #when treated samples have more than 1 value, can use comparison of 2 means
139         if (length(unique(treat$sampleID)) > 1) {
140             p.value=wilcox.test(none$value, treat$value)$p.value
141             output <- treat
142             output$p.val <- p.value
143
144             out <- rbind(out, output)
145         }
146
147         #when treated samples have only 1 value, do 1-sample t test, where treated
148         value
149         #is the true mean
150         else if (length(unique(treat$sampleID)) == 1) {
151             x.bar <- mean(none$value)
152             y.value <- treat$value
153             x.sd <- sd(none$value)
154             n <- length(unique(none$sampleID))
155
156             t.stat <- (x.bar - y.value)/(x.sd/sqrt(n))
157             p.value <- 2*pt(-abs(t.stat), n-1)
158
159             output <- treat

```

```

155     output$p.val <- p.value
156
157     out <- rbind(out, output)
158   }
159
160   #when the gentreat is not in that donor
161   else if (length(unique(treat$sampleID)) < 1) {
162     out <- out
163   }
164 }
165 return(out)
166 }
167
168 mw.result <- ddply(stat, c('donor', 'variable'),
169                   mw_custom, c('Clindamycin', 'RePOOPulate',
170                               'Clindamycin+RePOOPulate', 'Norepinephrine',
171                               'Vancomycin'))
172
173 #getting difference between none and treatment means
174
175 diff_custom <- function(data, gentreat) {
176   none <- data[data$gentreat == 'None',]
177   out <- none
178   out$uni.treat <- NA
179   out$diff <- NA
180   out$norm.diff <- NA
181   out$per.diff <- NA
182
183   for(i in gentreat) {
184     treat <- data[data$gentreat == i,]
185     if(nrow(treat) != 0) {
186       diff <- treat
187       diff$uni.treat <- paste(unique(treat$gentreat), unique(treat$donor), sep='_')
188       diff$diff <- unique(treat$avg) - unique(none$avg)
189       diff$norm.diff <- (unique(treat$avg) - unique(none$avg)) / (unique(none$avg) +
190                                                                    unique(treat$avg)) / 2 * 100
191       diff$per.diff <- (unique(treat$avg) - unique(none$avg)) / unique(none$avg) * 100
192
193       out <- rbind(out, diff)
194     } else if(nrow(treat) == 0) out <- out
195   }
196   return(out)
197 }
198
199 diff <- ddply(mw.result, c('donor', 'variable'),
200             diff_custom, c('Clindamycin', 'RePOOPulate',
201                           'Clindamycin+RePOOPulate', 'Norepinephrine', 'Vancomycin'))
202
203 #removing none rows
204
205 diff <- diff[diff$gentreat != 'None',]
206
207 #giving artificial diff values where non significant differences
208 #are assigned percent difference of zero
209 diff$sig.perdiff <- diff$per.diff
210 diff$sig.perdiff[diff$p.val > 0.05] <- 0
211
212 #artificially turning compounds with hi sig.perdiff into NA to make
213 it a different colour when plotting
214 diff$sig.perdiff[diff$sig.perdiff > 100 & !is.na(diff$sig.perdiff)] <- NA
215
216 #zero compounds
217
218 #Getting compounds with zero concentration in none or treated, but not both

```

```

218 zero.cmpd <- unique(as.character(diff$variable[!is.na(diff$per.diff) &
219                                     (diff$per.diff == -100 | diff$per.
220                                         diff == 'Inf')]))
221
222 zero.data <- diff[diff$variable %in% zero.cmpd,]
223
224 bin.list <- defined_binning()
225
226 #setting working directory back
227 setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript B/
    Analysis')
228
229 zero.order <- list()
230 for(i in 1:length(bin.list)) {
231   temp <- zero.data[zero.data$variable %in% bin.list[[i]],]
232   temp <- unique(temp[,c('uni.treat', 'variable', 'diff')])
233   temp$variable <- factor(temp$variable, levels=rev(levels(bin.list[[i]])), ordered
    =TRUE)
234   temp$uni.treat <- factor(temp$uni.treat,
235                           levels=c('Norepinephrine_A', 'RePOOPulate_A',
236                                   'Clindamycin+RePOOPulate_A',
237                                   'Clindamycin_A', 'Clindamycin_B', 'Vancomycin_RP2
    '),
238                                   ordered=TRUE)
239   temp <- temp[order(temp$variable),]
240   temp <- temp[order(temp$uni.treat),]
241
242   zero.order[[i]] <- temp
243 }
244
245 zero.export <- do.call(rbind.data.frame, zero.order)
246 write.csv(file='../Results/zero.data.csv', zero.export)
247
248 #hi compounds


---


249 #getting compounds that have sig.perdiff > 100
250 hi.cmpd <- unique(as.character(diff$variable[diff$per.diff > 100 &
251                                             !is.na(diff$per.diff)]))
252
253 hi.data <- diff[diff$variable %in% hi.cmpd,]
254 hi.data <- hi.data[!hi.data$variable %in% zero.cmpd,]
255
256 hi.order <- list()
257 for(i in 1:length(bin.list)) {
258   temp <- hi.data[hi.data$variable %in% bin.list[[i]],]
259   temp <- unique(temp[,c('uni.treat', 'variable', 'per.diff')])
260   temp$variable <- factor(temp$variable, levels=rev(levels(bin.list[[i]])), ordered
    =TRUE)
261   temp$uni.treat <- factor(temp$uni.treat,
262                           levels=c('Norepinephrine_A', 'RePOOPulate_A',
263                                   'Clindamycin+RePOOPulate_A',
264                                   'Clindamycin_A', 'Clindamycin_B', 'Vancomycin_RP2
    '),
265                                   ordered=TRUE)
266   temp <- temp[order(temp$variable),]
267   temp <- temp[order(temp$uni.treat),]
268
269   hi.order[[i]] <- temp
270 }
271
272 for(i in 1:length(hi.order)){
273   hi.write <- hi.order[[i]]
274   write.csv(file=paste('../Results/hi.data', i, '.csv', sep=''), hi.write)
275 }
276 #density plots


---


277 #looking at density plot of mean concentrations

```

```

278 plot(density(stat.diff$avg))
279 plot(density(stat.zero$avg))
280 plot(density(stat.hi$avg))
281 plot(density(diff.data$sig.perdiff[!is.na(diff.data$sig.perdiff)]))
282
283 #Assigning diff data to bins as defined by none-treated samples


---


284
285 bin.none <- bin.list
286 bin.treat <- bin.function(diff, bin=c(20,10,1,0.1,0))
287 bin.treat <- do.call(c, bin.treat)
288
289 none.cmpd <- c()
290 treat.cmpd <- c()
291 for(i in 1:length(bin.none)) {
292   temp <- as.character(bin.none[[i]])
293   temp2 <- as.character(bin.treat[[i]])
294   none.cmpd <- c(none.cmpd, temp)
295   treat.cmpd <- c(treat.cmpd, temp2)
296 }
297
298 #check if any compounds in treated samples not in non-treated samples
299 # if no can just use none sample binning
300 if(any(!treat.cmpd %in% none.cmpd)==FALSE) bin.list <- bin.none
301
302 #if yes, then need to augment none sample binning to include those compounds
303 #in the right order
304 ####write when needed####
305
306 #Setting treatment order


---


307
308 diff$gentreat <- factor(diff$gentreat,
309   levels=c('Norepinephrine', 'RePOOPulate',
310     'Clindamycin+RePOOPulate',
311     'Clindamycin', 'Vancomycin'))
312
313 diff$uni.treat <- factor(diff$uni.treat,
314   levels=c('Norepinephrine_A', 'RePOOPulate_A',
315     'Clindamycin+RePOOPulate_A',
316     'Clindamycin_A', 'Clindamycin_B', 'Vancomycin_RP2')
317 )
318
319 diff$pad <- sprintf('%022s', diff$variable)
320 #diff$variable <- factor(diff$variable, levels=none.cmpd, ordered=TRUE)
321
322 pad.list <- list()
323 for(i in 1:length(bin.none)) {
324   pad.temp <- sprintf('%022s', as.character(bin.list[[i]]))
325   pad.list[[i]] <- pad.temp
326 }
327
328 pad.order <- c()
329 for(i in 1:length(pad.list)) {
330   temp <- as.character(pad.list[[i]])
331   pad.order <- c(pad.order, temp)
332 }
333
334 #hi compounds for table


---


335
336 #Compounds with percent change > 100%
337 hi.cmpd <- unique(diff$variable[diff$per.diff > 100 & !is.na(diff$per.diff)])
338
339 #diff data for those compounds
340 out <- unique(diff[diff$variable %in% hi.cmpd, c('variable', 'uni.treat', 'per.diff',
341   'p.val')])
342
343 #converting per.diff to zero if insignificant
344 out$per.diff[out$p.val > 0.05] <- 0

```

```

342 #transforming into casted form
343 out2 <- dcast(out, variable~uni.treat, value.var='per.diff')
344
345 #re-ordering columns to particular order
346 out3 <- out2[,c('variable', 'Norepinephrine_A', 'RePOOPulate_A', 'Clindamycin+
      RePOOPulate_A',
347               'Clindamycin_A', 'Clindamycin_B', 'Vancomycin_RP2')]
348
349 #rearranging compounds to bin order
350 out3$variable <- factor(out3$variable, levels=none.cmpd, ordered=TRUE)
351 out3 <- out3[order(out3$variable),]
352
353 #Adding compound classes
354 type.info <- unique(cmpd.info[cmpd.info$Compound %in% hi.cmpd,1:2])
355 type.info$Compound <- factor(type.info$Compound, levels=none.cmpd, ordered=TRUE)
356 type.info <- type.info[order(type.info$Compound),]
357
358 out3$class <- type.info$Type
359
360 path.info <- cmpd.info[cmpd.info$Compound %in% hi.cmpd,]
361
362 #export as csv
363 #write.csv(file='../Results/hi_percent_difference.csv', out3)
364 write.csv(file='../Results/path_info.csv', path.info)
365
366 #exporting hi and zero data

```

```

367 #exporting hi data and zero data to csv file
368 zero.data2 <- unique(zero.data[,c('variable', 'uni.treat', 'med', 'diff')])
369
370
371 #ordering zero.data2 so variables are in order of median values
372 order.zero2 <- zero.data2[order(zero.data2$med, decreasing=TRUE),]
373 order.zero2$variable <- factor(order.zero2$variable,
374                               levels=rev(unique(order.zero2$variable)),
375                               ordered=TRUE)
376 zero.ordered2 <- zero.data2
377 zero.ordered2$variable <- factor(zero.ordered2$variable,
378                               levels=levels(order.zero2$variable))
379 zero.ordered2$uni.treat <- factor(zero.ordered2$uni.treat,
380                                 levels=c('Norepinephrine_A', 'RePOOPulate_A',
381                                           'Clindamycin+RePOOPulate_A',
382                                           'Clindamycin_A', 'Clindamycin_B',
383                                           'Vancomycin_RP2'))
384
385 zero.out <- dcast(zero.data2, formula=variable~uni.treat, value.var='diff')
386
387 #write.csv(file='../Plots/zero_data.csv', zero.out)
388
389 hi.data2 <- unique(hi.data[,c('variable', 'uni.treat', 'med', 'sig.perdiff')])
390
391 #ordering hi.data2 so variables are in order of median values
392 order.hi2 <- hi.data2[order(hi.data2$med, decreasing=TRUE),]
393 order.hi2$variable <- factor(order.hi2$variable, levels=rev(unique(order.hi2$
      variable))),
394                               ordered=TRUE)
395 hi.ordered2 <- hi.data2
396 hi.ordered2$variable <- factor(hi.ordered2$variable, levels=levels(order.hi2$
      variable))
397 hi.ordered2$uni.treat <- factor(hi.ordered2$uni.treat,
398                               levels=c('Norepinephrine_A', 'RePOOPulate_A',
399                                           'Clindamycin+RePOOPulate_A',
400                                           'Clindamycin_A', 'Clindamycin_B',
401                                           'Vancomycin_RP2'))
402
403 hi.out <- dcast(hi.ordered2, formula=variable~uni.treat, value.var='sig.perdiff')

```

```

404 #write.csv(file='check hi data.csv', hi.data)
405 #write.csv(file='check zero data.csv', zero.data)
406 #write.csv(file='../Plots/hi data.csv', hi.out)
407
408 #binning hi data by sig.perdiff


---


409 #binning hi.data by sig.perdiff
410 hi.data$bin <- NA
411 hi.data[hi.data$sig.perdiff >= 1000 & !is.na(hi.data$sig.perdiff), 'bin'] <- 'bin1'
412
413 hi.data[hi.data$sig.perdiff >= 500 & hi.data$sig.perdiff < 1000 &
414 !is.na(hi.data$sig.perdiff), 'bin'] <- 'bin2'
415 hi.data[hi.data$sig.perdiff >= 100 & hi.data$sig.perdiff < 500 &
416 !is.na(hi.data$sig.perdiff), 'bin'] <- 'bin3'
417
418 #converting compounds to character class and also gettinng compounds in each bin
419 bin.cmpd <- list()
420 for(i in c('bin1','bin2','bin3')) {
421   temp <- hi.data[hi.data$bin == i & !is.na(hi.data$bin),]
422   bin.cmpd[[i]] <- as.character(unique(temp$variable))
423 }
424
425 #Getting only compounds found in one bin and not any other bins before it
426 common.cmpd <- c(bin.cmpd[[1]])
427 hi.list <- list(bin.cmpd[[1]])
428
429 for(i in 2:length(bin.cmpd)) {
430   hi.list[[i]] <- bin.cmpd[[i]][!bin.cmpd[[i]] %in% common.cmpd]
431   common.cmpd <- c(common.cmpd, hi.list[[i]])
432 }
433
434 #ordering hi.data so variables are in order of median values
435 order.hi <- hi.data[order(hi.data$med, decreasing=TRUE),]
436 order.hi$variable <- factor(order.hi$variable, levels=rev(unique(order.hi$variable)
437 ),
438                               ordered=TRUE)
439
440 hi.ordered <- hi.data
441 hi.ordered$variable <- factor(hi.ordered$variable, levels=levels(order.hi$variable)
442 )
443
444 hi.ordered$gentreat <- factor(hi.ordered$gentreat,
445                               levels=levels(ordered.data$gentreat))
446 hi.ordered$uni.treat <- factor(hi.ordered$uni.treat,
447                               levels=levels(ordered.data$uni.treat))
448
449 #plotting


---


450 plot.list <- list()
451 legend.list <- list()
452 fig.list <- list()
453 fig.leg <- list()
454 leg.lim <- list()
455
456 leg.width <- c(10,10,10,10,10)
457 leg.pos <- list(c(0.5, 0.5), c(0.5, 0.5), c(0.5, 0.5), c(0.5, 0.5), c(0.6, 0.45))
458 fig.width <- c(8.6, 8.5, 8.6, 9.2, 9.4)
459 fig.height <- c(2.4, 2.8, 4.2, 13.5, 11.2)
460
461 margin <- list(c(0.1, 0.8, 0.9, 0.8),
462               c(0.1, 0.8, 1, 1),
463               c(0.1, 0.8, 0.9, 0.5),
464               c(0, 0.2, 0, 0),
465               c(0, 0.2, 0, 0))
466
467 #Panels 1:3


---


468 for(i in 1:3) {
469   plot.data <- unique(diff[diff$variable %in% bin.list[[i]],
470                         c('variable', 'pad', 'gentreat', 'uni.treat', 'p.val',
471                           'avg', 'med', 'diff', 'per.diff', 'sig.perdiff')])

```

```

468
469 #Setting compound order
470 plot.data$variable <- factor(as.character(plot.data$variable),
471                             levels=rev(none.cmpd[none.cmpd %in% bin.list[[i]]]),
472                             ordered=TRUE)
473 plot.data$pad <- factor(as.character(plot.data$pad),
474                         levels=rev(pad.order[pad.order %in% pad.list[[i]]]),
475                         ordered=TRUE)
476
477 plot.data$label <- '*'
478
479 #write.csv(file=paste('../Results/check_plot_data',i,'.csv',sep=''), temp)
480 #getting legend limits
481 max.lim <- max(plot.data$sig.perdiff[!is.na(plot.data$sig.perdiff)])
482 min.lim <- min(plot.data$sig.perdiff[!is.na(plot.data$sig.perdiff)])
483 larger.lim <- max(abs(max.lim), abs(min.lim))
484 leg.lim[[i]] <- c(-larger.lim, larger.lim)
485 print(max.lim)
486 print(min.lim)
487
488 plot.list[[i]] <- ggplot(plot.data, aes(x=uni.treat, y=pad, label=label))
489
490 plot.list[[i]] <- plot.list[[i]] + geom_tile(fill=NA, colour='white')
491
492 if(any(is.na(plot.data$sig.perdiff))==FALSE) {
493   plot.list[[i]] <- plot.list[[i]] +
494     geom_tile(aes(fill=sig.perdiff), colour='white')
495   print('here!')
496 }
497
498 else if(any(is.na(plot.data$sig.perdiff))==TRUE) {
499   plot.list[[i]] <- plot.list[[i]] +
500
501     #fill tile by sig.perdiff
502     geom_tile(data=plot.data[!is.na(plot.data$sig.perdiff),],
503              aes(fill=sig.perdiff), colour='white')
504
505   if(nrow(plot.data[plot.data$p.val <=0.05 &
506                    plot.data$per.diff == -100 &
507                    !is.na(plot.data$per.diff),]) > 0) {
508
509     plot.list[[i]] <- plot.list[[i]] +
510
511     #compounds with percent diff of -100% coloured with red with *
512     geom_tile(data=plot.data[plot.data$p.val <=0.05 &
513                             plot.data$per.diff == -100 &
514                             !is.na(plot.data$per.diff),],
515              colour='white', fill='red') +
516     geom_text(data=plot.data[plot.data$p.val <=0.05 &
517                             plot.data$per.diff == -100 &
518                             !is.na(plot.data$per.diff),],
519              colour='white', size=6, hjust=0.5, vjust=0.8)
520   }
521
522   if(nrow(plot.data[(plot.data$p.val <= 0.05 &
523                    plot.data$per.diff > 100 & !is.na(plot.data$sig.perdiff))
524                    |
525                    plot.data$per.diff == 'Inf',]) > 0) {
526
527     plot.list[[i]] <- plot.list[[i]] +
528
529     #compounds percent diff > 100% or Inf (0 in none, 1 in treated) coloured
530     #blue with *
531     geom_tile(data=plot.data[(plot.data$p.val <= 0.05 &
532                             plot.data$per.diff > 100 & !is.na(plot.data$
533                             sig.perdiff)) |
534                             plot.data$per.diff == 'Inf',],

```



```

532         colour='white', fill='blue') +
533     geom_text(data=plot.data[(plot.data$p.val <= 0.05 &
534         plot.data$per.diff > 100 & !is.nan(plot.data$
535             sig.perdiff)) |
536         plot.data$per.diff == 'Inf'],,
537         colour='white', size=6, hjust=0.5, vjust=0.8)
538     }
539     print('here2')
540 }
541 }
542 #-----
543 plot.list[[i]] <- plot.list[[i]] +
544     scale_x_discrete(labels=c('Norepinephrine', 'MET-1', 'Clindamycin+\nMET-1',
545         'Clindamycin (A)', 'Clindamycin (B)', 'Vancomycin')) +
546     scale_fill_gradientn(colours=c('red', 'white', 'blue'),
547         values=rescale(c(-larger.lim, 0, larger.lim)),
548         name='Percent Difference\nin Concentration\t',
549         limits=leg.lim[[i]]) +
550     theme_bw(12) +
551     theme(panel.background=element_rect('grey60'),
552         panel.grid.major=element_line(colour = 'grey40'),
553         axis.text.x=element_blank(),
554         axis.ticks.x=element_blank(),
555         axis.title.x=element_blank(),
556         axis.title.y=element_blank(),
557         legend.key.width=unit(1.1, 'cm'),
558         legend.key.height=unit(0.35, 'cm'),
559         legend.position=c(0.55, 0.4), legend.direction='horizontal',
560         legend.text=element_text(size=7.5),
561         legend.title=element_text(size=7.5),
562         plot.margin=unit(margin[[i]], 'cm'))
563
564 legend.list[[i]] <- g_legend(plot.list[[i]])
565
566 #Export as panels-----
567 if(FALSE){
568     #export as EPS; panels with legend
569     setEPS()
570     postscript(file=paste('../Results/test_plot', i, '.eps', sep=''),
571         onefile=TRUE, paper='special', horizontal=FALSE,
572         height=fig.height[i]+1, width=fig.width[i], pointsize=10)
573     grid.arrange(legend.list[[i]], plot.list[[i]] + theme(legend.position='none'),
574         heights=c(0.1, 0.9))
575     dev.off()
576
577     #export as EPS; panels separate from legend
578     setEPS()
579     postscript(file=paste('../Results/test_plot', i, '.eps', sep=''),
580         onefile=FALSE, paper='special', horizontal=FALSE,
581         height=fig.height[i], width=fig.width[i], pointsize=10)
582     fig.list[[i]] <- grid.arrange(plot.list[[i]] + theme(legend.position='none'))
583     dev.off()
584
585     #export as PNG - panels
586     png(filename=paste('../Plots/treatment effects/per diff_bin', i, '.png', sep=''),
587         type='cairo', height=fig.height[i], width=fig.width[i], units='cm', res
588             =240, pointsize=10)
589     fig.list[[i]] <- grid.arrange(plot.list[[i]] + theme(legend.position='none'))
590     dev.off()
591
592     png(filename=paste('../Plots/treatment effects/per diff_leg', i, '.png', sep=''),
593         height=1, width=leg.width[i], units='cm', res=240, pointsize=10)
594     fig.leg[[i]] <- grid.arrange(legend.list[[i]])
595     dev.off()
596 }
597 }

```

```

598 #panels 4 and 5
599 for(i in 4:length(bin.list)) {
600   plot.data <- unique(diff[diff$variable %in% bin.list[[i]],
601     c('variable','pad','gentreat','uni.treat','p.val',
602       'avg','med','diff','per.diff','sig.perdiff')])
603
604   #Setting compound order
605   plot.data$variable <- factor(as.character(plot.data$variable),
606     levels=rev(none.cmpd[none.cmpd %in% bin.list[[i]]]),
607     ordered=TRUE)
608   plot.data$pad <- factor(as.character(plot.data$pad),
609     levels=rev(pad.order[pad.order %in% pad.list[[i]]]),
610     ordered=TRUE)
611
612   plot.data$label <- '*'
613
614   #write.csv(file=paste('../Results/check_plot_data',i,'.csv',sep=''), temp)
615   #getting legend limits
616   max.lim <- max(plot.data$sig.perdiff[!is.na(plot.data$sig.perdiff)])
617   min.lim <- min(plot.data$sig.perdiff[!is.na(plot.data$sig.perdiff)])
618   larger.lim <- max(abs(max.lim), abs(min.lim))
619   leg.lim[[i]] <- c(-larger.lim, larger.lim)
620   print(max.lim)
621   print(min.lim)
622
623   plot.list[[i]] <- ggplot(plot.data, aes(x=uni.treat, y=pad, label=label))
624
625   plot.list[[i]] <- plot.list[[i]] + geom_tile(fill=NA, colour='white')
626
627   if(any(is.na(plot.data$sig.perdiff))==FALSE) {
628     plot.list[[i]] <- plot.list[[i]] +
629       geom_tile(aes(fill=sig.perdiff), colour='white')
630     print('here1')
631   }
632
633   else if(any(is.na(plot.data$sig.perdiff))==TRUE) {
634     plot.list[[i]] <- plot.list[[i]] +
635
636       #fill tile by sig.perdiff
637       geom_tile(data=plot.data[!is.na(plot.data$sig.perdiff),],
638         aes(fill=sig.perdiff), colour='white')
639
640     if(nrow(plot.data[plot.data$p.val <=0.05 &
641       plot.data$per.diff == -100 &
642       !is.na(plot.data$per.diff),]) > 0) {
643
644       plot.list[[i]] <- plot.list[[i]] +
645
646         #compounds with percent diff of -100% coloured with red with *
647         geom_tile(data=plot.data[plot.data$p.val <=0.05 &
648           plot.data$per.diff == -100 &
649           !is.na(plot.data$per.diff),],
650           colour='white', fill='red') +
651         geom_text(data=plot.data[plot.data$p.val <=0.05 &
652           plot.data$per.diff == -100 &
653           !is.na(plot.data$per.diff),],
654           colour='white', size=6, hjust=0.5, vjust=0.8)
655       }
656
657     if(nrow(plot.data[(plot.data$p.val <= 0.05 &
658       plot.data$per.diff > 100 & !is.na(plot.data$sig.perdiff)) |
659       plot.data$per.diff == 'Inf',]) > 0) {
660
661       plot.list[[i]] <- plot.list[[i]] +
662

```

```

663     #compounds percent diff > 100% or Inf (0 in none, 1 in treated) coloured
        blue with *
664     geom_tile(data=plot.data[(plot.data$p.val <= 0.05 &
665                             plot.data$per.diff > 100 & !is.nan(plot.data$
                                sig.perdiff)) |
                                plot.data$per.diff == 'Inf',],
666             colour='white', fill='blue') +
667     geom_text(data=plot.data[(plot.data$p.val <= 0.05 &
668                             plot.data$per.diff > 100 & !is.nan(plot.data$
                                sig.perdiff)) |
                                plot.data$per.diff == 'Inf',],
669             colour='white', size=6, hjust=0.5, vjust=0.8)
670 }
671
672 }
673
674 print('here2')
675 }
676 }
677 #

```

```

678 plot.list[[i]] <- plot.list[[i]] +
679     scale_x_discrete(labels=c('Norepinephrine', 'MET-1', 'Clindamycin+\nMET-1',
680                             'Clindamycin (A)', 'Clindamycin (B)', 'Vancomycin'))
681     +
682     scale_fill_gradientn(colours=c('red', 'white', 'blue'),
683                          values=rescale(c(-larger.lim, 0, larger.lim)),
684                          name='Percent Difference\nin Concentration\t',
685                          limits=leg.lim[[i]]) +
686     theme_bw(12) +
687     xlab('Treatment') +
688     theme(panel.background=element_rect('grey60'),
689           panel.grid.major=element_line(colour = 'grey40'),
690           axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
691           axis.title.y=element_blank(),
692           plot.margin=unit(c(0.5, 0.8, 0, 0), 'cm'),
693           legend.key.width=unit(1.1, 'cm'),
694           legend.key.height=unit(0.35, 'cm'),
695           legend.position=c(0.55, 0.4), legend.direction='horizontal',
696           legend.text=element_text(size=7.5),
697           legend.title=element_text(size=7.5),
698           plot.margin=unit(margin[[i]], 'cm'))
699
700 legend.list[[i]] <- g_legend(plot.list[[i]])
701
702 #Export as panels-----
703 if(FALSE){
704     #export as EPS; panels with legend
705     setEPS()
706     postscript(file=paste('../Results/test plot', i, '.eps', sep=''),
707               onefile=TRUE, paper='special', horizontal=FALSE,
708               height=fig.height[i]+1, width=fig.width[i], pointsize=10)
709     grid.arrange(legend.list[[i]], plot.list[[i]] + theme(legend.position='none')
710                 ,
711                 heights=c(0.1, 0.9))
712     dev.off()
713
714     #export as EPS; panels separate from legend
715     setEPS()
716     postscript(file=paste('../Results/test plot', i, '.eps', sep=''),
717               onefile=FALSE, paper='special', horizontal=FALSE,
718               height=fig.height[i], width=fig.width[i], pointsize=10)
719     fig.list[[i]] <- grid.arrange(plot.list[[i]] + theme(legend.position='none'))
720     dev.off()
721
722     #export as PNG - panels
723     png(filename=paste('../Plots/treatment effects/per diff_bin', i, '.png', sep=''),
724         type='cairo', height=fig.height[i], width=fig.width[i], units='cm', res

```

```

723         =240, pointsize=10)
724     fig.list[[i]] <- grid.arrange(plot.list[[i]] + theme(legend.position='none'))
725     dev.off()
726     png(filename=paste('.. /Plots/treatment effects/per diff_leg',i, '.png', sep=''),
727           height=1, width=leg.width[i], units='cm', res=240, pointsize=10)
728     fig.leg[[i]] <- grid.arrange(legend.list[[i]])
729     dev.off()
730 }
731 }
732
733
734 #Setting plot dimensions
735 widths <- list(c(10,6)/2.54+1,      #1
736               c(10,6)/2.54+1,      #2
737               c(10,6)/2.54+1,      #3
738               c(10,6)/2.54+1,      #4
739               c(10,10.4)/2.54+1)   #5
740
741 heights <- list(c(1,2.4)/2.54,      #1
742                c(1,3)/2.54,         #2
743                c(1,4.8)/2.54,       #3
744                c(1,15)/2.54,        #4
745                c(1,11.3)/2.54)      #5
746
747
748
749 #PNG
750 png(filename='.. /Plots/treatment effects/Percent Difference4.png',
751       width=10.5/2.54*2, height=sum(c(heights[[4]], heights[[3]])),
752       units='in', res=240, pointsize=10)
753 grid.arrange(arrangeGrob(legend.list[[1]],
754                           plot.list[[1]] + theme(legend.position='none'),
755                           legend.list[[2]],
756                           plot.list[[2]] + theme(legend.position='none'),
757                           legend.list[[5]],
758                           plot.list[[5]] + theme(legend.position='none'),
759                           ncol=1, heights=c(heights[[1]], heights[[2]], heights[[5]]),
760                           widths=c(widths[[1]], widths[[2]], widths[[5]])),
761             arrangeGrob(legend.list[[3]],
762                           plot.list[[3]] + theme(legend.position='none'),
763                           legend.list[[4]],
764                           plot.list[[4]] + theme(legend.position='none'),
765                           ncol=1, heights=c(heights[[3]], heights[[4]]),
766                           widths=c(widths[[3]], widths[[4]])),
767             ncol=2)
768 grid.text('A', x=unit(0.5, 'cm'), y=unit(21.1, 'cm'), gp=gpar(cex=1.5))
769 grid.text('C', x=unit(0.5, 'cm'), y=unit(17.4, 'cm'), gp=gpar(cex=1.5))
770 grid.text('E', x=unit(0.5, 'cm'), y=unit(12.95, 'cm'), gp=gpar(cex=1.5))
771 grid.text('B', x=unit(11, 'cm'), y=unit(21.2, 'cm'), gp=gpar(cex=1.5))
772 grid.text('D', x=unit(11, 'cm'), y=unit(15.4, 'cm'), gp=gpar(cex=1.5))
773
774 dev.off()
775
776 #EPS
777
778 setEPS()
779 postscript(file='.. /Results/test_plot.eps',
780           onefile=FALSE, paper='special', horizontal=FALSE,
781           width=10.5/2.54*2, height=sum(c(heights[[4]], heights[[3]])),
782           pointsize=10)
783 grid.arrange(arrangeGrob(legend.list[[1]],
784                           plot.list[[1]] + theme(legend.position='none'),
785                           legend.list[[2]],
786                           plot.list[[2]] + theme(legend.position='none'),
787                           legend.list[[5]],

```

```

788         plot.list[[5]] + theme(legend.position='none'),
789         ncol=1, heights=c(heights[[1]], heights[[2]], heights[[5]]),
790         widths=c(widths[[1]], widths[[2]], widths[[5]])),
791     arrangeGrob(legend.list[[3]],
792                 plot.list[[3]] + theme(legend.position='none'),
793                 legend.list[[4]],
794                 plot.list[[4]] + theme(legend.position='none'),
795                 ncol=1, heights=c(heights[[3]], heights[[4]]),
796                 widths=c(widths[[3]], widths[[4]])),
797     ncol=2)
798 dev.off()
799
800
801
802 #-----
803 #plotting percent diff where per.diff is > 100%
804
805 plot.list3 <- list()
806 leg.list3 <- list()
807 leg.lim3 <- list()
808
809 leg.pos3 <- list(c(0.5, 0.5), c(0.5, 0.5), c(0.5, 0.5), c(0.5, 0.5), c(0.6, 0.45))
810 leg.width3 <- c(8,8,8,9,10)
811
812 fig.width3 <- c(8.7, 8.5, 9.3, 9.8, 11)
813 fig.height3 <- c(4.3,3.9,8.2,17.5,13.5)
814 fig.prop3 <- list(c(0.3, 0.7), c(0.4, 0.96), c(0.05, 0.95), c(0.05, 0.95), c
815                   (0.05,0.95))
816 for(i in 1:length(hi.list)) {
817     plot.data <- hi.ordered[hi.ordered$variable %in% hi.list[[i]],
818                             c('variable','gentreat','uni.treat','p.val',
819                               'avg','med','diff','per.diff','sig.perdiff')]
820
821     #write.csv(file=paste('check hi data plotted',i,'.csv',sep=''), plot.data)
822
823     #getting legend limits
824     hi.max <- max(plot.data$sig.perdiff[!is.na(plot.data$sig.perdiff)])
825     hi.min <- min(plot.data$sig.perdiff[!is.na(plot.data$sig.perdiff)])
826     hi.larger <- max(hi.max, hi.min)
827     leg.lim3[[i]] <- c(-hi.larger, hi.larger)
828
829     plot.list3[[i]] <- ggplot(plot.data, aes(x=uni.treat, y=variable))
830
831     if(any(is.nan(plot.data$sig.perdiff))==FALSE) {
832         plot.list3[[i]] <- plot.list3[[i]] +
833             geom_tile(aes(fill=sig.perdiff), colour='white')
834         print('here1')
835     }
836
837     else if(any(is.nan(plot.data$sig.perdiff))==TRUE) {
838         plot.list3[[i]] <- plot.list3[[i]] +
839             geom_tile(data=plot.data[!is.na(plot.data$sig.perdiff),],
840                       aes(fill=sig.perdiff), colour='white') +
841             geom_tile(data=plot.data[is.nan(plot.data$sig.perdiff),], colour='white',
842                       fill='grey30', alpha=0.5)
843         print('here2')
844     }
845
846     plot.list3[[i]] <- plot.list3[[i]] +
847         scale_x_discrete(labels=c('Norepinephrine','Repopulate','Clindamycin+\n\
848                                   tRepopulate',
849                                   'Clindamycin (A)','Clindamycin (B)','Vancomycin')) +
850         scale_fill_gradientn(colours=c('red','white','blue'),
851                              values=rescale(c(-hi.larger,0,hi.larger)),
852                              name='Percent Difference\nin Concentration\t',
853                              limits=leg.lim3[[i]] +
854

```

```

854   xlab('Treatment') +
855   ylab('Compound') +
856   theme(panel.grid.major=element_line(colour = 'grey40'),
857         axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
858         legend.position=leg.pos3[[i]], legend.direction='horizontal',
859         legend.key.width=unit(1, 'cm'),
860         legend.key.height=unit(0.3, 'cm'),
861         legend.text=element_text(size=7.5),
862         legend.title=element_text(size=7.5))
863
864   leg.list3[[i]] <- g_legend(plot.list3[[i]])
865
866   png(filename=paste('../Plots/treatment effects/per diff-hi range',i,'.png',sep=''),
867        height=fig.height3[i], width=fig.width3[i], units='cm', res=240, pointsize
868        =10)
869   fig <- grid.arrange(plot.list3[[i]] + theme(legend.position='none'))
870   dev.off()
871   png(filename=paste('../Plots/treatment effects/per diff-hi_leg',i,'.png',sep=''),
872        height=1, width=leg.width3[i], units='cm', res=240, pointsize=10)
873   leg <- grid.arrange(leg.list3[[i]])
874   dev.off()
875 }
876
877 #-----
878 plot.list2 <- list()
879 legend.list2 <- list()
880 fig.list2 <- list()
881 fig.leg2 <- list()
882 leg.lim2 <- list()
883 plot.data.list <- list()
884
885 #customizing plot specifications
886 leg.pos2 <- list(c(0.5, 0.5), c(0.5, 0.5), c(0.5, 0.5), c(0.5, 0.5), c(0.6, 0.45))
887 leg.width2 <- c(8,8,8,9,10)
888
889 fig.width2 <- c(8.7, 8.5, 9.3, 8, 11)
890 fig.height2 <- c(4.3,3.9,8.2,8,13.5)
891 fig.prop2 <- list(c(0.3, 0.7), c(0.4, 0.96), c(0.05, 0.95), c(0.05, 0.95), c
892                  (0.05,0.95))
893
894 for(i in 1:(length(uni.list2)-1)) {
895   plot.data <- unique(ordered.data2[ordered.data2$variable %in% uni.list2[[i]],])
896
897   plot.data.list[[i]] <- plot.data
898
899   #getting legend limits
900   max.lim <- max(plot.data$diff[!is.na(plot.data$diff)])
901   min.lim <- min(plot.data$diff[!is.na(plot.data$diff)])
902   larger.lim <- max(abs(max.lim), abs(min.lim))
903   leg.lim2[[i]] <- c(-larger.lim, larger.lim)
904   print(max.lim)
905   print(min.lim)
906   plot.list2[[i]] <- ggplot(plot.data, aes(x=uni.treat, y=variable))
907
908   if(any(is.na(plot.data$diff))==FALSE) {
909     plot.list2[[i]] <- plot.list2[[i]] +
910       geom_tile(aes(fill=diff), colour='white')
911   }
912
913   else if(any(is.na(plot.data$diff))==TRUE) {
914     plot.list2[[i]] <- plot.list2[[i]] +
915       geom_tile(data=plot.data[!is.na(plot.data$diff),],
916                aes(fill=diff), colour='white') +
917       geom_tile(data=plot.data[is.na(plot.data$diff),], colour='white',
918                fill='grey30', alpha=0.5)
919   }

```

```

919 }
920 plot.list2[[i]] <- plot.list2[[i]] +
921   scale_x_discrete(labels=c('Norepinephrine', 'Repopulate', 'Clindamycin+\n\
    tRepopulate',
922     'Clindamycin (A)', 'Clindamycin (B)', 'Vancomycin')) +
923   scale_fill_gradientn(colours=c('red', 'white', 'blue'),
924     values=rescale(c(-larger.lim, 0, larger.lim)),
925     name='Difference\nin Concentration\t',
926     limits=leg.lim2[[i]]) +
927   theme_bw(10) +
928   xlab('Treatment') +
929   ylab('Compound') +
930   theme(panel.grid.major=element_line(colour = 'grey40'),
931     axis.text.x=element_text(hjust=0, vjust=1, angle=-45),
932     legend.position=leg.pos2[[i]], legend.direction='horizontal',
933     legend.key.width=unit(1, 'cm'),
934     legend.key.height=unit(0.3, 'cm'),
935     legend.text=element_text(size=7.5),
936     legend.title=element_text(size=7.5))
937
938 legend.list2[[i]] <- g_legend(plot.list2[[i]])
939
940 png(filename=paste('../Plots/treatment effects/diff_bin', i, '.png', sep=''),
941   height=fig.height2[i], width=fig.width2[i], units='cm', res=240, pointsize
    =10)
942 fig.list2[[i]] <- grid.arrange(plot.list2[[i]] + theme(legend.position='none'))
943 dev.off()
944
945 png(filename=paste('../Plots/treatment effects/diff_leg', i, '.png', sep=''),
946   height=1, width=leg.width2[i], units='cm', res=240, pointsize=10)
947 fig.leg2[[i]] <- grid.arrange(legend.list2[[i]], heights=fig.prop2[[i]])
948 dev.off()
949
950 }

```

C.6 Functions

function-summary preprocessing.R

```

1  #####
2  #Function for preprocessing exported files that contain
3  #molecular weight of compounds and pH of samples
4
5  #input: file - is the csv file name, along with its path
6  #         object - is the object name you want to store the processed raw data in
7  #output: data - data frame of the concentration summaries with appropriate column
    names
8  #         and row names
9  #         mol.wt - molecular weights of compounds
10
11 summary.preprocess <- function(file.name){
12
13   #reading csv file
14   data <- read.csv(file=file.name, header=FALSE, check.names=FALSE,
15     encoding='UTF-8')
16
17   #removing compound molecular weights and sample pH values
18   mol.data <- data[c(3,4), c(1,3:ncol(data))]
19
20   col.name <- as.character(as.matrix((mol.data[1,2:ncol(mol.data)])))
21   row.name <- as.character(as.matrix(mol.data[2,1]))
22
23   mol.weight <- mol.data[2,2:ncol(mol.data)]
24   colnames(mol.weight) <- col.name
25   rownames(mol.weight) <- row.name

```

```

26
27 data <- data[c(3,5:nrow(data)),c(1,3:ncol(data))]]
28
29 #creating vector of column names
30 col.name <- as.matrix(cbind('file', data[1,2:ncol(data)]))
31 colnames(data) <- col.name
32
33 #removing redundant row of compounds, compounds are now column names
34 data <- data[2:nrow(data),]
35
36 #turning all concentrations into numeric class
37 for(i in 2:ncol(data)) data[,i] <- as.numeric(as.character(data[,i]))
38
39 #turning all file names into characters
40 data$file <- as.character(data$file)
41
42 return(list('data'=data, 'mol.wt'=mol.weight))
43
44 }

```

function-data cleanup.R

```

1 #####
2 #This function-Removing contaminants and zero-columns
3 data_cleanup = function(dataset, sampleID='sampleID', variable='variable',
4 value='value', form.y, form.x, contaminants=NULL,
5 export=FALSE, convert.NA=FALSE) {
6
7 #loading required packages
8 require(reshape2)
9 require(plyr)
10 #Getting list of compounds before anything is removed
11 pre = unique(dataset[,variable])
12
13 #Removing DSS
14 dataset = dataset[!dataset[,variable] %in% c('DSS-d6 (Chemical Shape Indicator)',
15 'DSS (Chemical Shape Indicator)'),]
16
17 check.dss = any(dataset[,variable] %in% c('DSS-d6 (Chemical Shape Indicator)',
18 'DSS (Chemical Shape Indicator)'))
19
20 #Removing contaminants
21 dataset = dataset[!dataset[,variable] %in% contaminants,]
22
23 #Getting list of compounds after DSS and contaminants removed
24 post = unique(dataset[,variable])
25
26 #Checking if contaminants were removed
27 check.contaminants = pre[!pre %in% post]
28 check.contaminants = check.contaminants[!check.contaminants %in% c('DSS-d6 (
29 Chemical Shape Indicator)',
30 'DSS (Chemical Shape Indicator)')]
31
32 #Removing compounds with zero concentrations in every sample
33 subset = ddply(dataset, variable, summarise, sample=sampleID,
34 value=value)
35
36 zero.values <- ddply(subset, variable, summarise,
37 unique.values=length(unique(value)))
38
39 constant <- zero.values[zero.values$unique.values == 1,]
40
41 #Removing those rows
42 dataset = dataset[!dataset[,variable] %in% constant[,variable],]

```



```

43 if(export==TRUE) {
44   #Building the formula needed to cast data to wide format
45   form.y <- paste(form.y, collapse=' + ')
46   form = as.formula(paste(form.y, form.x, sep=' ~ '))
47
48   wide <- dcast(data=dataset, formula=form, value.var=value)
49   check2 <- dcast(data=dataset, formula=form, value.var=value)
50
51   write.csv(wide, file='datacheck_wide.csv')
52   write.csv(dataset, file='datacheck_long.csv')
53 }
54
55 if(convert.NA==TRUE) {
56   #converting NAs into 0
57   dataset$value[is.na(dataset$value)] <- 0
58   write.csv(dataset, file='clean data check.csv')
59 }
60
61 #Printing out confirmations
62 cat('Compounds with zero or constant concentration in all samples:\n')
63 cat(paste('\t', constant[,variable], collapse='\n'),'\n')
64 cat('Are there Compounds with na concentration?\n')
65 cat(any(is.na(dataset[,value])),'\n')
66
67
68
69 cat('DSS included?\n')
70 cat(check.dss,'\n')
71 cat('Contaminants removed:\n')
72 cat(paste('\t', check.contaminants, collapse='\n'),'\n')
73 return(dataset)
74 }

```

function-descriptive statistics.R

```

1 #####
2 #Function to obtain descriptive data of data
3 # calculates mean, standard deviation, standard error, ymax, ymin
4 # has option to scale data, and corresponding descriptive stats
5
6 #input: full.data - qualifier and concentration data
7 #       data.type - 'long' or 'wide'; specifies format of full.data
8 #       variable - character string of how you want long data to be broken up;
9 #               not required for wide data
10 #       qualifier - qualifier data; used so can separate out from full.data
11 #       scale.data - logical; retrieve stats on the scaled version of the data;
12 #               mean centres and unit variance scales the data
13
14 #output: stat - data frame of descriptive statistics
15 #----->
16 descriptive.stat <- function(full.data, data.type='wide', variable=NULL,
17                             qualifier=NULL, scale.data=FALSE) {
18
19   require(plyr)
20   if(data.type=='wide') {
21     #separating concentration data from qualifier data in full.data
22     data <- full.data[, (ncol(qualifier)+1):ncol(full.data)]
23
24     if(scale.data==FALSE) {
25       stat <- data.frame(avg = apply(data, 2, mean),
26                           med = apply(data, 2, median),
27                           stdev = apply(data, 2, sd),
28                           compound = colnames(data))
29
30       stat$se <- stat$stdev/sqrt(nrow(data))
31       stat$ymax <- stat$avg + stat$se
32       stat$ymin <- stat$avg - stat$se
33     }

```

```

34
35     else if(scale.data==TRUE){
36         stat <- data.frame(scale.avg = apply(data, 2, mean),
37                             scale.med = apply(data, 2, med),
38                             scale.stdev = apply(data, 2, sd),
39                             compound = colnames(data))
40
41         stat$scale.se <- stat$scale.stdev/sqrt(nrow(data))
42         stat$scale.ymax <- stat$scale.avg + stat$scale.se
43         stat$scale.ymin <- stat$scale.avg - stat$scale.se
44     }
45 }
46
47 else if(data.type=='long') {
48     if(scale.data==FALSE) {
49         stat <- ddply(full.data, variable, mutate,
50                       avg = mean(value),
51                       med = median(value[value != 0]),
52                       stdev = sd(value),
53                       n=length(value))
54
55         stat$se <- stat$stdev/sqrt(stat$n)
56         stat$ymax <- stat$avg + stat$se
57         stat$ymin <- stat$avg - stat$se
58
59         #turning all NAs to zero in median
60         stat$med[is.na(stat$med)] <- 0
61
62     }
63     else if(scale.data==TRUE) {
64         stat <- ddply(full.data, variable, mutate,
65                       scale.avg = mean(value),
66                       scale.med = median(value[value != 0]),
67                       scale.stdev = sd(value))
68
69         stat$scale.se <- stat$scale.stdev/sqrt(nrow(full.data))
70         stat$scale.ymax <- stat$scale.avg + stat$scale.se
71         stat$scale.ymin <- stat$scale.avg - stat$scale.se
72     }
73 }
74 }
75 return(stat)
76 }

```

function-extract legend.R

```

1 #####
2 #function to extract legend from a ggplot
3 #to be used the grid.arrange from gridExtra package
4
5 #Extract Legend
6 g_legend<-function(a.gplot){
7   tmp <- ggplot_gtable(ggplot_build(a.gplot))
8   leg <- which(sapply(tmp$grobs, function(x) x$name) == "guide-box")
9   legend <- tmp$grobs[[leg]]
10  return(legend)}

```

function-pca.R

```

1 #####
2
3 #Functions for making processing and analyzing data subsets created earlier
4
5 #----->
6 #Function for Performing PCA analysis
7 #when running the script:

```

```

8  #enter as: pca_processed_data —> process to be applied to the output_numeric
   data, PCA
9
10
11 #input: data – numeric data to be used for PCA, called output_numeric, unpackaged
   from 'function-data input' file
12 #   form.y – y variable in formula; usually is qualifier column names to be
   put into formula that casts data
13 #   form.x – x variable in formula; usually the variables involved (i.e.
   compound)
14 #   variable – the column name in which variables are stored (i.e. 'Compound
   ')
15 #   value – the column name in which data values are stored (i.e. '
   concentration')
16 #   scale – 'UV' for unit variance, 'UVN' for no unit variance, 'pareto'
17 #   center – mean centering; default to TRUE
18 #output: PCA-processed data, called pca_data
19 #----->
20 #Loading required packages
21 require(reshape2)
22
23 perform_pca = function(data, sampleID='sample_name', variable, value,
24                        form.y, form.x, scale, center = TRUE, convert.NA=FALSE) {
25
26
27
28 #Getting column names for qualifiers and for concentrations
29 numeric.long = c(variable, value)
30 qual.long = colnames(data)[!colnames(data) %in% numeric.long]
31
32 #Building the formula needed to cast data to wide format
33 form.y <- paste(form.y, collapse=' + ')
34 form = as.formula(paste(form.y, form.x, sep=' ~ '))
35
36 #Casting data to wide format; includes qualifiers
37 data.wide = dcast(data=data, formula=form,
38                  value.var=value)
39
40 if(convert.NA==TRUE) {
41   #converting NAs into 0
42   data.wide[is.na(data.wide)] <- 0
43   write.csv(data.wide, file='data_check2.csv')
44 }
45
46 #separating out concentration and qualifiers in wide format
47 conc.wide = data.wide[, unique(data[, variable])]
48 qual.wide = data.wide[, qual.long]
49
50 #Converting concentration data to numeric class
51 for(i in conc.wide) { i <- as.numeric(i)}
52
53 #Unit Variance Scaling, scaling weight of 1/sk
54 if(scale == 'UV') {
55   scale = TRUE
56 }
57 #No Unit Variance Scaling
58 else if(scale == 'UVN') {
59   scale = FALSE
60 }
61 #Pareto scaling, scaling weight of 1/(sk)e-2
62 else if(scale == 'pareto') {
63   scale = for(column in 1:ncol(data)) {
64     data[, column] = data[, column]/sqrt(sd(data[, column]))
65   }
66 }
67
68 #Creating pca model
69 pca_data = prcomp(conc.wide, scale=scale, center=center,

```

```

70      na.action=na.omit)
71
72      print(summary(pca_data))
73
74      #Adding qualifiers to pca score data
75      pca_data$x = cbind(qual.wide, pca_data$x)
76
77      #Melting pca score data back to long format,
78      #where variable is PC number and value is score value on that PC
79      pca.melt = melt(pca_data$x, id.vars=qual.long)
80
81      #Removing 'PC' in PC column names, so only left with a number
82      pca.melt$variable = as.numeric(gsub('PC', '', pca.melt$variable))
83
84      #Get rid of last PC if it is odd
85      if(max(pca.melt$variable) %% 2 == 1) {
86          pca.melt = pca.melt[pca.melt$variable != max(pca.melt$variable), ]
87      }
88
89      #Creating a new object where has columns of of odd number PCs and even number PCs
90      odd.PCs = pca.melt$variable %% 2 == 1
91      split = pca.melt[odd.PCs, qual.long]
92      split$x.value = pca.melt$value[odd.PCs]
93      split$y.value = pca.melt$value[!odd.PCs]
94      #Creating a column called PC, to give the PC number
95      split$x.PC = pca.melt$variable[odd.PCs]
96      split$y.PC = pca.melt$variable[!odd.PCs]
97      #Adding column 'view' to describe which PCs are being plotted
98      split$view = paste('PC',pca.melt$variable[odd.PCs], 'vs',
99                        'PC',pca.melt$variable[!odd.PCs])
100
101      #Replacing the pca model score data with the melted and sorted version created
102      #here
103      pca_data$x = split
104
105      return(pca_data)
106  }

```

function-pca diagnostics.R

```

1  #####
2
3  #Functions for performing diagnostic tests on the pca data, evaluating which PCs to
4  #consider
5  #----->
6
7  #Function for performing PCA diagnostics
8  #when running the script:
9  #enter as = diagnostic_results
10 #unpack as diagnostic_results_number = diagnostics_data[['number']]
11 # diagnostics_results_explained = diagnostics_data[['explained']]
12 # diagnostics_results_plot = diagnostics_data[['plot']]
13
14 #input: the PCA model output of prcomp()
15 #Output: varies, defined by individual function
16
17 #----->
18 #Scree plot diagnostic
19 #
20 #Output: 'number' - the number of PCA components to use
21 # 'explained' - % variance explained by 'number' PCs
22 # 'plot' - ggplot object with generated scree plot
23
24 scree_plot = function(model)
25 {
26     #Making sure the right packages are loaded
27     require(ggplot2)

```

```

28
29 #Testing if the model is indeed the result of a prcomp function
30 if (class(model) != 'prcomp')
31 {
32   stop('Input to scree-plot must be the output of prcomp() function')
33 }
34
35 #Converting standard deviations to variance
36 variance = model$sdev**2
37 variance = variance/sum(variance)*100
38
39 #Looking at differences between sequential variance values
40 difference = variance[-length(variance)] - variance[2:length(variance)]
41
42 #Creating a permanent dataset that will be used later for plotting
43 #for now, every value is assigned INCLUDE, and INCLUDE is defined as TRUE
44 dataset = data.frame(x=1:length(difference),
45                     y=difference,
46                     INCLUDE=TRUE)
47
48 #Looping through data to exclude PC variance differences, and creating a straight
49 line fit
50 #with every iteration
51 for (i in 1:min(10, nrow(dataset)-2))
52 {
53   #Marking point for exclusion, making all INCLUDE values equal to FALSE
54   dataset[i, 'INCLUDE'] = FALSE
55
56   #Calculating length of x vector (in sum, TRUE is equal to 1, FALSE, 0)
57   x_length = sum(dataset$INCLUDE)
58
59   #Generating weights for linear regression. We want the first points to have a
60   lot
61   #more weight than the last points. To do this, we create a weight vector from 1
62   to
63   #0.1 in equal segments, then cube it to amplify the effect.
64   weight = seq(1, .1, length.out=x_length)**3
65
66   #Fitting linear model, lm = to fit linear model, y ~ x means y=f(x), where y is
67   the difference in variance between two successive values
68   #fitting the linear model with the dataset assigned previously, and using all
69   rows with INCLUDE in its column
70   model = lm(y ~ x, data=dataset[dataset$INCLUDE,])
71
72   #Extracting probability that the slope is actually 0
73   prob = summary(model)$coefficients[2,4]
74
75   #see objects of the model by >print(summary(model)).
76   #Column 2 of object Coefficients is the Std. Error
77   #Column 4 of object Coefficients is Pr(>|t|), which in context, gives the
78   probability that the slope is zero
79
80   #If the probability is sufficiently high, the slope
81   #is deemed to be zero and there is limited improvement
82   #from subsequent PC components; exit script
83   if (prob > 0.05) break
84 }
85
86 #Storing the last PC index (i) before the loop terminated
87 n = i
88
89 #Calculating how much variance this number
90 #of components will explain
91 explained = sum(variance[1:n])
92
93 #Creating a new variance dataset from diagnostic results
94 diagnostic_results = data.frame(x=1:length(variance),

```

```

90         y=variance ,
91         Type=ifelse(c(dataset$INCLUDE, TRUE), #test if INCLUDE in
           the dataset is TRUE
           'Not significant', #if yes, type out 'Not
           significant',
           'Significant')) #if no, type out "Significant
92
93
94
95 #Creating a basic ggplot object as validation
96 p = ggplot(diagnostic_results)
97 p = p + geom_point(aes(x=x, y=y, colour=Type))
98 p = p + stat_smooth(data=diagnostic_results[diagnostic_results$Type == 'Not
           significant'],,
           aes(x=x, y=y), se=FALSE, method='lm', colour='black')
99
100
101 p = p + xlab('Principal Component')
102 p = p + ylab('Variance explained (%)')
103 p = p + scale_colour_manual(values=c('Significant'='red', 'Not significant'='
           black'))
104
105 return(list('number'=n, 'explained'=explained, 'plot'=p))
106 }

```

function-plotting D score.R

```

1 #####
2
3 #Functions for plotting pca results
4
5 #List of defined functions:
6 # score_plot_2d
7
8 #Making sure the right packages are loaded
9 require(ggplot2)
10 require(grid)
11 require(plyr)
12 require(reshape2)
13 require(ellipse)
14
15 #sourcing required functions
16 source('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Functions/
           function-plotting 1D score.R')
17
18 #----->
19
20 #Function for plotting 2D score plots
21 #when running the script:
22 #enter as: score_plot
23
24 #Input: the PCA model output of prcomp()
25 #Output: A single ggplot2 plot object
26
27 #Note: Here, the use of "2D" implies that 2 principal components
28 #are considered at a time. As the plotting is actually
29 #quite simple, this is no more than a shortcut function for
30 #repeated use. Two different modes are used. In 'full' mode,
31 #each principal component is compared against each other i.e.
32 #PC 1 vs PC 2, PC 1 vs PC 3, etc... In 'compressed' mode, only
33 #successive principal components are compared i.e. PC 1 vs PC 2,
34 #PC 3 vs PC 4, etc...
35 #----->
36
37 #Input: 'pca_data' - the PCA model output of prcomp()
38 # 'PCs' - vector of Principal Components to consider
39 # 'colour' - name of qualifer_data column to be used as the colour aesthetic
40 # 'shape' - name of qualifer_data column to be used as the shape aesthetic
41 # 'size' - name of qualifer_data column to be used as the size aesthetic
42 # 'label' - name of qualifer_data column to be used as the point labels

```

```

43     # 'view_mode' - the view at which the principal components are compared;
44     #           either compressed or full,
45     #           where compressed is PC1 v PC2, PC3 v PC4 etc, and full is PC1 v
46     #           PC2, PC1 v PC3, PC1 v PC4 etc.
47     # 'black-white' - plot will be in black and white
48 #Output: 'plot' - ggplot object with generated plot
49
50 #----->
51 score_plot_2d = function(pca_data, PCs=1:4,
52                          colour=NULL, shape=NULL, size=NULL, label=NULL, title=NULL
53                          ,
54                          legend.ori=NULL, outline=FALSE, black.white=FALSE)
55 {
56     #----->
57     #Input testing
58
59     #Testing if the pca_data is indeed the result of a prcomp function
60     if (class(pca_data) != 'prcomp')
61     {
62         stop('Input "pca_data" to score_plot_2d must be the output of prcomp() function
63             ')
64     }
65
66     #Testing qualifier rownames to make sure they are the same as the rownames of the
67     #score matrix
68     if (length(intersect(rownames(pca_data$x), rownames(qualifier_data))) != length(
69         rownames(qualifier_data)))
70     {
71         stop('Input "qualifier_data" to score_plot_2d must have the same rownames as
72             the score matrix in "pca_data"')
73     }
74
75     #Testing if PCs are numeric
76     if (!is.numeric(PCs))
77     {
78         stop('Input "PCs" to score_plot_2d must be a vector of numeric values
79             corresponding to target principal components')
80     }
81
82     #Testing to make sure that PCs are in the input data
83     if (any(PCs > ncol(pca_data$x)))
84     {
85         stop('Input "PCs" to score_plot_2d must not be outside of the pca_data
86             principal component range')
87     }
88
89     #Checking if colour, shape, or size conditions are column names in qualifer_data
90     aesthetics = c('colour'=colour, 'shape'=shape, 'size'=size)
91     aesthetics = aesthetics[!is.null(aesthetics)] #Excluding default NULL entries
92     #----->
93
94     #Extract scores
95     score.data = as.data.frame(pca_data$x)
96
97     #Copying out the standard deviation info from pca_data as a dataframe
98     variance = as.data.frame(pca_data$sdev)
99     #Converting standard deviations to variance
100     variance = variance**2
101     variance = variance/sum(variance)*100
102     variance = round(variance, digits=1)
103
104     #Initializing a new vector to store collected.data
105     collected.data = c()
106
107     #The way in which the collected data is built up and
108     #If only 1 PC was specified, calling
109     #score_plot_1d instead

```

```

102 if (length(PCs) == 1)
103 {
104   p = score_plot_1d(pca_data, PCs,
105                     colour, shape, size, label, pt_lab, title,
106                     perform_checks=FALSE, group_labels)
107   return(p)
108 }
109
110 #Iterating through the PCs, two at a time.
111 for (i in seq(1, length(PCs), by=2))
112 {
113   view.data = score.data[score.data$x.PC==i,]
114   view.data$view.var = paste('PC', i, '(', variance[i], '%) ', 'vs. ',
115                              'PC', i+1, '(', variance[i+1], '%)')
116
117   #Adding to collected.data
118   collected.data = rbind(collected.data, view.data)
119 }
120
121 #Converting View to factor, preserving ordering
122 collected.data$view.var = factor(collected.data$view.var,
123                                 levels=unique(collected.data$view.var))
124
125 #Calculating plot limits
126 xlimit = max(abs(collected.data$x.value))*1.15
127 ylimit = max(abs(collected.data$y.value))*1.1
128
129 #
130 #Generating plot
131
132
133 p = ggplot(collected.data)
134
135 #If size descriptor has been specified, use it as an aesthetic,
136 if(!is.null(size)) {
137   p = p + geom_point(aes_string(x='x.value', y='y.value',
138                                 colour=colour, shape=shape, size=size), alpha=0.8)
139 }
140 else
141 { #otherwise, default to 3 and make points with black outline
142   if(outline==TRUE){
143     p = p + geom_point(aes_string(x='x.value', y='y.value', shape=shape),
144                       colour='black', size=4, alpha=0.8)
145   }
146
147   p = p + geom_point(aes_string(x='x.value', y='y.value', colour=colour, shape=
148     shape),
149                     size=3, alpha=0.8)
150   #Alternative script for making points with black outline
151   #though there is internal bug with legend
152   #p = p + geom_point(aes_string(x='xPC', y='yPC', fill=colour, shape=shape),
153   #
154   #p = p + scale_shape_manual(values=21:25)
155 }
156
157 #Cusomizing colours of points
158 #if colour is discrete give customized colours
159 if(!is.numeric(collected.data[, colour]) & length(unique(collected.data[, colour]))
160 < 8) {
161   p = p + scale_colour_manual(values=c('deeppink1', 'darkorange1', 'cornflowerblue',
162   'darkorchid4', 'chartreuse', 'blue', 'grey30'))
163 }
164 else if(!is.numeric(collected.data[, colour]) & length(unique(collected.data[,
165   colour])) >= 8) {
166   p = p + scale_colour_brewer(palette='Dark2')

```



```

165 }
166 }
167 #if colour is continuous give colour gradient
168 else if (is.numeric(collected.data[,colour])) {
169
170   p = p + scale_colour_gradient(low='red')
171 }
172
173
174 #Adding text only if the "label" descriptor has been provided
175 if (!is.null(label))
176 {
177   if (label == 'Plot_Treat') {
178     p = p + geom_text(aes_string(x='x.value', y='y.value', label=label, alpha='
179       Plot_Treat'),
180                       hjust=-0.05, vjust=0.75,
181                       size=5, fontface='plain', lineheight=0.8)
182     p = p + scale_alpha_manual(values=list('None'=0, 'Norepinephrine'=1,
183       'Clindamycin'=1, 'Repopulate'=1,
184       'Clindamycin+RePOPulation'=1), guide=
185         FALSE)
186   }
187   else {
188     p = p + geom_text(aes_string(x='x.value', y='y.value', label=label),
189                       hjust=-0.5, vjust=0.3,
190                       size=5, fontface='plain', lineheight=0.8)
191   }
192 }
193
194 #Creating scheme for black and white plots
195 if (black_white == TRUE){
196   p = p + geom_point(aes_string(x='x.value', y='y.value', shape=shape, fill=
197     colour),
198                     size=3, alpha=0.8) +
199   scale_shape_manual(values=21:25) +
200   scale_fill_manual(values=c('grey30', 'white'))
201 }
202
203 #Adding title only if "title" descriptor has been provided
204 if (!is.null(title)) p = p + ggtitle(label=title)
205
206 #Drawing origin
207 p = p + geom_vline(xintercept=0, alpha=0.3)
208 p = p + geom_hline(yintercept=0, alpha=0.3)
209
210 p = p + xlab('PC score') +
211 ylab('PC score') +
212 theme_bw(16) +
213 theme(axis.title.x = element_text(size=14, face='bold'),
214       axis.title.y = element_text(size=14, face='bold'),
215       axis.text = element_text(size=12, colour='black'),
216       legend.title = element_text(size=14),
217       legend.text = element_text(size=14),
218       legend.key = element_rect(colour='white'),
219       legend.margin = unit(0, "cm"),
220       legend.position = 'bottom',
221       panel.border=element_rect(size=2, colour='black'),
222       panel.margin=unit(.05, 'npc'))
223
224 #Legend orientation
225 if (legend.ori == 'horizontal') {
226   p = p + theme(legend.direction='horizontal',
227     legend.position='bottom')
228 }
229 else if(legend.ori == 'vertical') {
230   p = p + theme(legend.direction='vertical',
231     legend.position='right')
232 }

```

```

230 }
231
232 #Setting limits
233 p = p + xlim(-xlim, xlim)
234 p = p + ylim(-ylim, ylim)
235
236 #While the rest of the plotting code is generic,
237 #the facetting will change based on view
238 #if (view_mode == 'compressed')
239 p = p + facet_wrap(~ view.var, ncol=min(floor(length(PCs)/2), 3)) +
240 theme(strip.text.x = element_text(size=12, face='bold'))
241 #else if (view_mode == 'full') p = p + facet_grid(yLabel ~ xLabel)
242
243 return(list('collected.data'=collected.data, 'p'=p))
244 }

```

function-plotting D loading.R

```

1 #####
2 #Function for plotting2D Loading plot
3 #
4 #Here, the use of "2D" implies that 2 principal components
5 #are considered at a time. Unlike the score plot, loading plot groupings
6 #(pathways) are not exclusive i.e. a compound may belong to more than one
7 #grouping. This makes it difficult to use colour, shape, and other aesthetics
8 #in the same manner as for the score plot. For now, these options will be
9 #removed
10 #
11 #when running the script:
12 #enter as: loading_plot
13 #
14 #Input: 'pca_data' - the PCA model output of prcomp()
15 # 'PCs' - vector of Principal Components to consider
16 # 'label' - vector of compounds that should be labelled on the score plot
17 # 'view_mode' - one of either 'full' or 'compressed' (see score plot)
18 #
19 #Output: 'plot' - ggplot object with generated plot
20
21 loading_plot_2d = function(pca_data, PCs=1:4, label=NULL, outline=FALSE,
22                             title=NULL, view_mode='compressed')
23 {
24   #----->
25   #Input testing
26
27   #Testing if the pca_data is indeed the result of
28   #a prcomp function
29   if (class(pca_data) != 'prcomp')
30   {
31     stop('Input "pca_data" to score_plot_2d must be the output of prcomp() function')
32   }
33
34   #Testing if PCs are numeric
35   if (!is.numeric(PCs))
36   {
37     stop('Input "PCs" to score_plot_2d must be a vector of numeric values
38           corresponding to target principal components')
39   }
40
41   #Testing to make sure that PCs are in the input data
42   if (any(PCs > ncol(pca_data$rotation)))
43   {
44     stop('Input "PCs" to score_plot_2d must not be outside of the pca_data
45           principal component range')
46   }
47
48   #Checking if given labels are in the loading data
49   if (any(!label %in% rownames(pca_data$rotation)))

```

```

48 {
49   stop(paste('The following labels do not correspond to variables in "pca_data":\n\t',
50             paste(label[!label %in% rownames(pca_data$rotation)], collapse=', ')
51             ))
52 }
53 #----->
54 #Extract loadings
55 loading_data = as.data.frame(pca_data$rotation)
56
57 variance = as.data.frame(pca_data$sdev)
58 #Converting standard deviations to variance
59 variance = variance**2
60 variance = variance/sum(variance)*100
61 variance = round(variance, digits=1)
62
63 #Initializing a new vector to store collected_data
64 collected_data = c()
65
66 #The way in which the collected data is built up and
67 #presented depends on the view_model
68 if (view_mode == 'compressed')
69 {
70
71   #If only 1 PC was specified, calling
72   #loading_plot_1d instead
73   if (length(PCs) == 1)
74   {
75     p = loading_plot_1d(pca_data, PCs, label, title,
76                        perform_checks=FALSE)
77     return(p)
78   }
79   #Iterating through the PCs, two at a time.
80   #If there is an odd number of PCs, the last one is dropped
81   for (i in seq(1, length(PCs)-1, by=2))
82   {
83     view_data = loading_data[, c(i, i+1)]
84     view_data$View = paste('PC', i, '(', variance[i], '%)', 'vs.',
85                           'PC', i+1, '(', variance[i+1], '%)')
86
87     #Changing score_column names
88     colnames(view_data)[1:2] = c('xPC', 'yPC')
89
90     #Initializing text column with blanks
91     view_data$Text = ''
92
93     #Adding labels
94     view_data[label, 'Text'] = label
95
96     #Adding to collected_data
97     collected_data = rbind(collected_data, view_data)
98   }
99 }
100 else if (view_mode == 'full')
101 {
102   #Iterating through the PCs one at a time. As the comparison
103   #is factorial, two loops are used, 1v2, 1v3, ..., 2v1, 2v2, etc
104   for (i in PCs)
105   {
106     for (j in PCs)
107     {
108       view_data = loading_data[, c(i, j)]
109       view_data$xLabel = paste('PC', i, '(', variance[i], '%)')
110       view_data$yLabel = paste('PC', j, '(', variance[j], '%)')
111
112       #Changing score_column names
113       colnames(view_data)[1:2] = c('xPC', 'yPC')

```

```

114
115     #Initializing text column with blanks
116     view_data$Text = ''
117
118     #Adding labels
119     view_data[label, 'Text'] = label
120
121     #Adding to collected_data
122     collected_data = rbind(collected_data, view_data)
123   }
124 }
125 }
126 else
127 {
128   stop('Input "view_mode" must be one of "full" or "compressed"')
129 }
130
131 #Calculating plot limits
132 xlimit = max(abs(collected_data$xPC))*1.1
133 ylimit = max(abs(collected_data$yPC))*1.1
134
135 #Generating plot
136 p = ggplot(collected_data)
137
138 #If highlighting will be performed, use a less stark, grey tone,
139 #otherwise, use black with transparency to limit overplotting
140 if (is.null(label)) p = p + geom_point(aes(x=xPC, y=yPC), colour='black', size=5,
141   alpha=.7)
142 else if (is.null(label) & outline == TRUE) {
143   p = p + geom_point(aes(x=xPC, y=yPC), colour='grey', size=5)
144 }
145 #Adding text only if the "label" descriptor has been provided,
146 #the points with text are also highlighted
147 if (!is.null(label))
148 {
149   p = p + geom_text(aes(x=xPC, y=yPC, label=Text),
150     hjust=-0.1, vjust=0, size=5, fontface=2)
151 }
152 if (!is.null(label) & outline==TRUE) {
153   p = p + geom_point(data=collected_data[collected_data$Text != '',],
154     aes(x=xPC, y=yPC), size=5, shape=1)
155 }
156
157 #Adding title only if "title" descriptor has been provided
158 if (!is.null(title)) p = p + ggtitle(label=title)
159
160 #Drawing origin
161 p = p + geom_vline(xintercept=0, alpha=0.3)
162 p = p + geom_hline(yintercept=0, alpha=0.3)
163
164 p = p + xlab('PC loading') + ylab('PC loading')
165
166 #Setting limits
167 p = p + xlim(-xlimit, xlimit)
168 p = p + ylim(-ylimit, ylimit)
169
170 #While the rest of the plotting code is generic,
171 #the facetting will change based on view
172 if (view_mode == 'compressed') p = p + facet_wrap(~ View, ncol=min(floor(length(
173   PCs)/2), 3)) +
174   theme_bw(16) +
175   theme(axis.title.x = element_text(size=14, face='bold'),
176     axis.title.y = element_text(size=14, face='bold'),
177     axis.text = element_text(size=12, colour='black'),
178     legend.title = element_text(size=14),
179     legend.text = element_text(size=12),
180     legend.key = element_rect(colour='white'),

```

```

180     legend.margin = unit(0, "cm"),
181     panel.border=element_rect(size=2, colour='black'),
182     strip.text.x = element_text(size=14, face='bold'))
183 else if (view_mode == 'full') p = p + facet_grid(yLabel ~ xLabel) +
184     theme_bw(16) +
185     theme(axis.title.x = element_text(size=14, face='bold'),
186           axis.title.y = element_text(size=14, face='bold'),
187           axis.text = element_text(size=12, colour='black'),
188           legend.title = element_text(size=14),
189           legend.text = element_text(size=12),
190           legend.key = element_rect(colour='white'),
191           legend.margin = unit(0, "cm"),
192           strip.text.x = element_text(size=14, face='bold'))
193
194 return(list('p'=p, 'collected.data'=collected_data))
195 }

```

function-scaling.R

```

1 #####
2 #This function is to scale metabolite data before it goes to PCA processing,
3 #or without PCA processing:
4 #   Scaling to Acetate
5 #   Mean centering
6 #   Unit variance scaling
7 # Note: mean centering and unit variance is done through the perform_pca function,
8 #       so if passing data from
9
10 #input: data - unprocessed concentration data
11 #        normalize - scale all compounds within one observation to a single compound
12 #               divides every compound by the concentration of normalizing
13 #               compound
14 #output: scaledData - concentration data that has undergone the specified scaling
15 #
16 #####
17 /
18 require(reshape2)
19
20
21
22 scaling = function(data, sampleID='sampleID', variable='variable', value='value',
23                   normalize=NULL, centre=TRUE, UV=FALSE) {
24
25   if(!is.null(normalize)){
26     setwd('C:/Users/Sandi/Dropbox/Liquid Gold/R Scripts/V 3.0/')
27     source('./Functions/function-normalize.R')
28     normalized = ddply(data, sampleID, f_normalize, variable=variable, scal=value,
29                       normalize=normalize)
30
31     #remove the compound by which the concentrations are scaled, otherwise PCA will
32     #fail
33     processed.data = normalized[normalized[,variable] != normalize, ]
34     print('data is normalized by Acetate')
35     print('Acetate removed from analysis as its value is one')
36   }
37   if(is.null(normalize)){
38     processed.data = data
39   }
40
41   if(centre==TRUE){
42     processed.data = ddply(processed.data, variable, mutate,
43                           scal=value-mean(value))
44     print('data is mean centred')
45   }
46 }

```

```

45  if(centre==FALSE){
46    processed.data = processed.data
47  }
48
49  if(UV==TRUE){
50    processed.data = ddply(processed.data, variable , mutate,
51                          scal=scal/sd(scal))
52    print('data is scaled by unit variance')
53  }
54  if(UV==FALSE){
55    processed.data = processed.data
56  }
57
58  return(processed.data)
59 }

```

function-t test.R

```

1  #####
2  #####
3  #function to perform the t test
4
5  #First getting the t statistic for each compound
6  tStatistic = function(data, value, variable, group1, group2, n) {
7
8    #Separating out population and test groups
9    pop = data[data[, variable] == group1,]
10   pop = pop[,colnames(pop) == value]
11
12   sample = data[data[, variable] == group2,]
13   sample = sample[,colnames(sample) == value]
14
15   #Calculating t statistic
16   pop = ddply(pop, 'compound', transform, average = mean(concentration))
17   pop = ddply(pop, 'compound', mutate, stdev = sd(concentration))
18   pop = pop[,colnames(pop)%in% c('compound', 'average', 'stdev')]
19   pop = unique(pop)
20   print(head(pop))
21
22   if(FALSE) {
23     sample = ddply(sample, 'compound', transform, average = mean(concentration))
24     sample = sample[,colnames(sample) %in% c('compound', 'average')]
25
26     calc1 = data.frame(compound = pop1.calc$compound,
27                        sub = abs(pop1.calc$average-pop2.calc$average))
28     calc2 = data.frame(compound = control.calc$compound,
29                        ss = control.calc$stdev*sqrt(2/n))
30     ss = data.frame(compound = pop1.calc$compound, ss = calc1$sub/calc2$ss)
31   }
32
33   return(list('pop'=pop, 'sample'=sample,
34              'calc1'=calc1, 'calc2'=calc2, 'ss'=ss))
35 }
36
37 #abs(apply(control,2,mean)-apply(treatment,2,mean)/(apply(control, 2, sd)*sqrt(2/
38   n)))
39
40 if(FALSE) {
41   test = tStatistic(comp.input, value='composition', variable='Gen-Treat',
42                    group1='None', group2='Clindamycin', 2)
43
44   control.calc = test[['pop1.calc']]
45   treat.calc = test[['pop2.calc']]
46   calc1 = test[['calc1']]
47   calc2 = test[['calc2']]
48   ss = test[['ss']]
49 }
50
51 if(FALSE){
52   t_test_func = function(data, value, variable, group1, group2) {

```

```

49 x = data[data[, variable] == group1,]
50 x = x[,colnames(x) == value]
51
52 y = data[data[, variable] == group2,]
53 y = y[,colnames(y) == value]
54
55 test = t.test(x, y)
56 out = test$p.value
57 return(out)
58 }
59 }

```

function-siegel tukey test.R

```

1 siegel.tukey <- function(x, y, id.col = FALSE, adjust.median = F,
2   rnd = -1, alternative = "two.sided", mu = 0, paired =
3   FALSE,
4   exact = FALSE, correct = TRUE, conf.int = FALSE, conf.
5   level = 0.95) {
6   ##### published on:
7   # http://www.r-statistics.com/2010/02/siegel-tukey-a-non-parametric-test-for-
8   # equality-in-variability-r-code/
9   ## Main author of the function: Daniel Malter
10
11  # x: a vector of data
12
13  # y: Group indicator (if id.col=TRUE); data of the second
14  # group (if
15  # id.col=FALSE). If y is the group indicator it MUST take 0
16  # or 1 to indicate
17  # the groups, and x must contain the data for both groups.
18
19  # id.col: If TRUE (default), then x is the data column and y
20  # is the ID column,
21  # indicating the groups. If FALSE, x and y are both data
22  # columns. id.col must
23  # be FALSE only if both data columns are of the same length.
24
25  # adjust.median: Should between-group differences in medians
26  # be leveled before
27  # performing the test? In certain cases, the Siegel-Tukey
28  # test is susceptible
29  # to median differences and may indicate significant
30  # differences in
31  # variability that, in reality, stem from differences in
32  # medians.
33
34  # rnd: Should the data be rounded and, if so, to which
35  # decimal? The default
36  # (-1) uses the data as is. Otherwise, rnd must be a
37  # non-negative integer.
38  # Typically, this option is not needed. However,
39  # occasionally, differences in
40  # the precision with which certain functions return values
41  # cause the merging
42  # of two data frames to fail within the siegel.tukey
43  # function. Only then
44  # rounding is necessary. This operation should not be
45  # performed if it affects
46  # the ranks of observations.
47
48  # ... arguments passed on to the Wilcoxon test. See
49  # ?wilcox.test
50
51  # Value: Among other output, the function returns the data,
52  # the Siegel-Tukey
53  # ranks, the associated Wilcoxon's W and the p-value for a
54  # Wilcoxon test on

```

```

52 # tie-adjusted Siegel-Tukey ranks (i.e., it performs and
53 # returns a
54 # Siegel-Tukey test). If significant, the group with the
55 # smaller rank sum has
56 # greater variability.
57
58 # References: Sidney Siegel and John Wilder Tukey (1960) "A
59 # nonparametric sum
60 # of ranks procedure for relative spread in unpaired
61 # samples." Journal of the
62 # American Statistical Association. See also, David J.
63 # Sheskin (2004)
64 # "Handbook of parametric and nonparametric statistical
65 # procedures." 3rd
66 # edition. Chapman and Hall/CRC. Boca Raton, FL.
67
68 # Notes: The Siegel-Tukey test has relatively low power and
69 # may, under certain
70 # conditions, indicate significance due to differences in
71 # medians rather than
72 # differences in variabilities (consider using the argument
73 # adjust.median).
74
75 # Output (in this order)
76
77 # 1. Group medians (after median adjustment if specified)
78 # 2. Wilcoxon-test for between-group differences in medians
79 # (after the median
80 # adjustment if specified)
81 # 3. Data, group membership, and the Siegel-Tukey ranks
82 # 4. Mean Siegel-Tukey rank by group (smaller values indicate
83 # greater
84 # variability)
85 # 5. Siegel-Tukey test (Wilcoxon test on tie-adjusted
86 # Siegel-Tukey ranks)
87
88 is.formula <- function(x) class(x) == "formula"
89
90 if (is.formula(x)) {
91   y <- do.call(c, list(as.name(all.vars(x)[2])), envir = parent.frame(2))
92   x <- do.call(c, list(as.name(all.vars(x)[1])), envir = parent.frame(2)) # I am
93   # using parent.frame(2) since if the name of the variable in the equation is
94   # 'x', then we will mistakenly get the function in here instead of the
95   # vector.
96   id.col <- TRUE
97   # print(x)
98   # print(ls.str())
99   # data=data.frame(c(x,y),rep(c(0,1),c(length(x),length(y))))
100  # print(data)
101 }
102
103 if (id.col == FALSE) {
104   data = data.frame(c(x, y), rep(c(0, 1), c(length(x), length(y))))
105 } else {
106   data = data.frame(x, y)
107 }
108 names(data) = c("x", "y")
109 data = data[order(data$x), ]
110 if (rnd > -1) {
111   data$x = round(data$x, rnd)
112 }
113
114 if (adjust.median == T) {
115   ##cat("\n", "Adjusting medians...", "\n", sep = "")
116   data$x[data$y == 0] = data$x[data$y == 0] - (median(data$x[data$y ==
117   0]))
118   data$x[data$y == 1] = data$x[data$y == 1] - (median(data$x[data$y ==
119   1]))

```



```

117 }
118 ##cat("\n", "Median of group 1 = ", median(data$x[data$y == 0]),
119 ##      "\n", sep = "")
120 ##cat("Median of group 2 = ", median(data$x[data$y == 1]), "\n",
121 ##      "\n", sep = "")
122 ##cat("Testing median differences...", "\n")
123 ##print(wilcox.test(data$x[data$y == 0], data$x[data$y == 1]))
124
125 # The following must be done for the case when id.col==F
126 x <- data$x
127 y <- data$y
128
129 ##cat("Performing Siegel-Tukey rank transformation...", "\n",
130 ##      "\n")
131
132
133
134 sort.x <- sort(data$x)
135 sort.id <- data$y[order(data$x)]
136
137 data.matrix <- data.frame(sort.x, sort.id)
138
139 base1 <- c(1, 4)
140 iterator1 <- matrix(seq(from = 1, to = length(x), by = 4)) -
141   1
142 rank1 <- apply(iterator1, 1, function(x) x + base1)
143
144 iterator2 <- matrix(seq(from = 2, to = length(x), by = 4))
145 base2 <- c(0, 1)
146 rank2 <- apply(iterator2, 1, function(x) x + base2)
147
148 #print(rank1)
149 #print(rank2)
150
151 if (length(rank1) == length(rank2)) {
152   rank <- c(rank1[1:floor(length(x)/2)], rev(rank2[1:ceiling(length(x)/2)]))
153 } else {
154   rank <- c(rank1[1:ceiling(length(x)/2)], rev(rank2[1:floor(length(x)/2)]))
155 }
156
157
158 unique.ranks <- tapply(rank, sort.x, mean)
159 unique.x <- as.numeric(as.character(names(unique.ranks)))
160
161 rank.matrix <- data.frame(unique.x, unique.ranks)
162
163 ST.matrix <- merge(data.matrix, rank.matrix, by.x = "sort.x",
164   by.y = "unique.x")
165
166 ##print(ST.matrix)
167
168 ##cat("\n", "Performing Siegel-Tukey test...", "\n", sep = "")
169
170 ranks0 <- ST.matrix$unique.ranks[ST.matrix$sort.id == 0]
171 ranks1 <- ST.matrix$unique.ranks[ST.matrix$sort.id == 1]
172
173 ##cat("\n", "Mean rank of group 0: ", mean(ranks0), "\n", sep = "")
174 ##cat("Mean rank of group 1: ", mean(ranks1), "\n", sep = "")
175
176 ##print(wilcox.test(ranks0, ranks1, alternative = alternative,
177 ##      mu = mu, paired = paired, exact = exact, correct = correct,
178 ##      conf.int = conf.int, conf.level = conf.level))
179
180 #putting function outputs as objects so can be extracted
181 mean.rank0 <- mean(ranks0)
182 mean.rank1 <- mean(ranks1)
183 st.test <- wilcox.test(ranks0, ranks1, alternative = alternative,
184   mu = mu, paired = paired, exact = exact, correct = correct

```

```

185         conf.int = conf.int , conf.level = conf.level)
186
187     return( list ( 'st.ranks' = ST.matrix , 'mean.rank0' = mean.rank0 ,
188                  'mean.rank1' = mean.rank1 , 'st.test' = st.test  ) )
189 }
190
191
192
193
194
195
196
197 if(F) {
198
199     #Example:
200
201     ### 1
202     x=c(4,4,5,5,6,6)
203     y=c(0,0,1,9,10,10)
204     siegel.tukey(x,y, F)
205     siegel.tukey(x,y) #same as above
206
207     ### 2
208     # example for a non equal number of cases:
209     x=c(4,4,5,5,6,6)
210     y=c(0,0,1,9,10)
211     siegel.tukey(x,y,F)
212
213     ### 3
214     x <- c(33, 62, 84, 85, 88, 93, 97, 4, 16, 48, 51, 66, 98)
215     id <- c(0,0,0,0,0,0,0,1,1,1,1,1,1)
216     siegel.tukey(x,id,T)
217     siegel.tukey(x~id) # from now on, this also works as a function...
218     siegel.tukey(x,id,T,adjust.median=F,exact=T)
219
220     ### 4
221     x<-c
222         (177,200,227,230,232,268,272,297,47,105,126,142,158,172,197,220,225,230,262,270)
223
224     id<-c(rep(0,8),rep(1,12))
225     siegel.tukey(x,id,T,adjust.median=T)
226
227     ### 5
228     x=c(33,62,84,85,88,93,97)
229     y=c(4,16,48,51,66,98)
230     siegel.tukey(x,y)
231
232     ### 6
233     x<-c(0,0,1,4,4,5,5,6,6,9,10,10)
234     id<-c(0,0,0,1,1,1,1,1,1,0,0,0)
235     siegel.tukey(x,id,T)
236
237     ### 7
238     x <- c(85,106,96, 105, 104, 108, 86)
239     id<-c(0,0,1,1,1,1,1)
240     siegel.tukey(x,id,T)
241 }

```

stat custom.R

```

1 #####
2 //
3 #Function to get descriptive statistics
4 #meant to be used inside a ddply

```

```

5 #getting mean normalized values for each compound
6 stat_custom <- function(d, value) {
7   stat <- data.frame(avg = mean(d[, value]),
8                       med = median(d[, value][d[, value] != 0]),
9                       stdev = sd(d[, value]),
10                      n=length(d[, value]))
11
12   stat$sse <- stat$stdev/sqrt(stat$n)
13   stat$ymax <- stat$avg + stat$sse
14   stat$ymin <- stat$avg - stat$sse
15
16   stat$r_max <- range(d[, value])[2]
17   stat$r_min <- range(d[, value])[1]
18   stat$per_r_max <- abs(stat$r_max-stat$avg)/stat$avg*100
19   stat$per_r_min <- abs(stat$r_min-stat$avg)/stat$avg*100
20
21   out <- cbind(d, stat)
22   return(out)
23 }

```

linear model.R

```

1 #####
2 /
3 #Test for linear model
4 #Extracts estimate, stanadard error, test statistic, p-value for slope and
5 intercept
6 lm_custom <- function(data, group, value='value'){
7   d <- data
8   f <- paste(value, group, sep='~')
9
10  result <- lm(as.formula(f), data=d)
11  result.sum <- summary(result)
12
13  out <- data.frame(slope=result.sum$coefficients[2,1],
14                   slope.se=result.sum$coefficients[2,2],
15                   slope.t=result.sum$coefficients[2,3],
16                   slope.p=result.sum$coefficients[2,4],
17                   int=result.sum$coefficients[1,1],
18                   int.se=result.sum$coefficients[1,2],
19                   int.t=result.sum$coefficients[1,3],
20                   int.p=result.sum$coefficients[1,4])
21
22  return(out)
23 }

```

normalize.R

```

1 #####
2 #Function to normalize concentrations by one compound
3 #written for Fecal water samples
4
5 norm_custom <- function(d.stat, variable='variable', value='value',
6                          norm='Acetate', norm.val='avg', norm.name='norm') {
7   norm.value <- unique(d.stat[, norm.val][d.stat[, variable]==norm])
8   d.stat[, norm.name] <- d.stat[, value]/norm.value
9   out <- unique(d.stat[, !colnames(d.stat) %in%
10                  c('med', 'stdev', 'n', 'se', 'ymax', 'ymin')])
11   return(out)
12 }

```

multiple linear regression.R

```

1 #####
2 #multiple linear regression

```

```

3 #for comparison of bias and interaction between more than one independent variables
4
5 mlr_custom <- function(d, f){
6   l.model <- lm(formula=f, data=d)
7   l.result <- summary(l.model)
8   l.coef <- as.data.frame(l.result$coefficient)
9   rname <- data.frame(rname=rownames(l.coef))
10  out <- cbind(rname, l.coef)
11  out$variable <- as.character(unique(d$variable))
12  return(out)
13 }

```

mlr prediction.R

```

1 #####
2 #Predicting data points along the plane of a multiple linear regression
3 #obtains mlr model first, then applies data points based on that model
4 mlr_predict <- function(d, f){
5   l.model <- lm(formula=f, data=d)
6   d$prediction <- predict(l.model)
7   return(d)
8 }

```

linear model.R

```

1 #####
2 /
3 #Test for linear model
4 #Extracts estimate, stanadard error, test statistic, p-value for slope and
5 intercept
6 lm_custom <- function(data, group){
7   d <- data
8   f <- paste('value', group, sep='~')
9
10  result <- lm(as.formula(f), data=d)
11  result.sum <- summary(result)
12
13  out <- data.frame(slope=result.sum$coefficients[2,1],
14                   slope.se=result.sum$coefficients[2,2],
15                   slope.t=result.sum$coefficients[2,3],
16                   slope.p=result.sum$coefficients[2,4],
17                   int=result.sum$coefficients[1,1],
18                   int.se=result.sum$coefficients[1,2],
19                   int.t=result.sum$coefficients[1,3],
20                   int.p=result.sum$coefficients[1,4],
21                   res.sum=sum(result.sum$residuals))
22
23  return(out)
24 }

```

ANOVA-mean.R

```

1 #####
2 //
3 #ANOVA
4 #Testing variation about the mean within and between groups
5 aov_custom <- function(data, group) {
6   d <- data
7
8   grand <- mean(d$value) #grand mean
9
10  ss.grp <- ddply(d, group, mutate, avg=mean(value)) #within-group mean
11  ss.grp$grand <- grand #adding grand mean
12 }

```

```

13  ssr <- sum((ss.grp$value - ss.grp$avg)**2) #sum of sqr of within groups
14  ssa <- sum((ss.grp$avg - ss.grp$grand)**2) #sum of sqr of between groups
15
16  dfr <- nrow(d)-length(unique(d[,group])) #n-q
17  dfa <- length(unique(d[,group]))-1 #q-1
18
19  msa <- ssa/dfa #mean square between group
20  msr <- ssr/dfr #mean square within group
21
22  f <- msa/msr #f statistic
23  p.val <- 1-pf(f, dfa, dfr)
24  out <- data.frame(variable=unique(d$variable), f.val=f, p.val=p.val)
25  return(out)
26 }

```

kruskal wallis.R

```

1  #####
2  /
3  #kruskal-Wallis for UC expt 2
4  #kruskal wallis test
5  kruskal_custom <- function(data, group){
6    d <- data
7    result <- kruskal.test(data$value, as.factor(data[,group]))$p.val
8    d$p.val <- result
9    return(d)
10 }

```

tukey.R

```

1  #####
2  //
3  #Tukey's HSD function for UC Expt 2
4
5  tukey_custom <- function(data, group) {
6    grps <- unique(data[,group])
7    pair <- combn(grps,2)
8    out <- c()
9    for(i in 1:ncol(pair)) {
10     grp1 <- data[data[,group] == pair[1,i],]
11     grp2 <- data[data[,group] == pair[2,i],]
12
13     grand <- mean(data$value) #grand mean
14     ss.grp <- ddply(data, group, mutate, avg=mean(value)) #within-group mean
15     ss.grp$grand <- grand #adding grand mean
16     ssr <- sum((ss.grp$value - ss.grp$avg)**2) #sum of sqr of within groups
17     dfr <- nrow(data)-length(unique(data[,group])) #n-q
18
19     avg1 <- mean(grp1$value)
20     avg2 <- mean(grp2$value)
21     n <- nrow(grp1)
22     se <- sqrt(ssr/n)
23
24     q.stat <- (max(c(avg1, avg2))-min(c(avg1, avg2)))/se
25     p.val <- ptukey(q.stat, n, df=dfr)
26
27     output <- data.frame(g1=pair[1,i], g2=pair[2,i], p.val=p.val)
28     out <- rbind(out, output)
29   }
30   return(out)
31 }

```

correlation.R

```

1  #####
2  //

```

```

2 #Correlation relationship
3
4 #Returns correlation coefficient , p-value and test statistic
5
6
7 corr_custom <- function(data, x, y, methods='pearson') {
8   d <- data
9   x <- as.numeric(data[,x])
10  y <- as.numeric(data[,y])
11  corr.result <- cor.test(x, y, method=methods, alternative='two.sided')
12  p.val <- corr.result$p.value
13  stat <- corr.result$statistic
14  estimate <- corr.result$estimate
15  out <- data.frame(stat=stat, estimate=estimate, p.val=p.val)
16  return(out)
17 }

```

Student t test.R

```

1 #####
2 //
3 #Function to perform Student's t test
4
5 t_custom <- function(data, group) {
6
7   grps <- unique(data[,group])
8   grp1 <- data[data[,group] == grps[1],]
9   grp2 <- data[data[,group] == grps[2],]
10
11  result <- t.test(grp1$value, grp2$value)
12  t.stat <- result$statistic
13  p.val <- result$p.value
14
15  out <- c(variable=unique(data$variable), t.stat=t.stat, p.val=p.val)
16
17  return(out)
18 }

```

function-heat map bin.R

```

1 #####
2 #####
3 #Function to divide data into bins, and have them returned in a particular order
4 # best to be used for plotting heat maps
5 # only makes 5 bins
6
7 #input: data = dataframe that contains the value to be represented in plot; long
8 form
9 # bin = vector of 5 values are the bin limits, in decreasing order
10 # variable = name of column that holds the variables to be represented in the
11 plot
12 # order.data = dataframe that holds the values by which the compounds should
13 be ordered;
14 usually is the same as data
15 # order.value = name of column that holds values by which order.var should be
16 ordered
17 # order.var = name of column the variables to be ordered;
18 # order.var and variable should hold the same info
19
20 #output: uni.list = class list; list of compounds in each bin
21 # ordered.data = dataframe of data to be plotted, and compounds are ordered
22 # by a decreasing value, usually mean concentration
23
24 bin_heatmap <- function(data, bin, variable, order.data, order.var, order.value) {

```

```

23 #Assigning data to bins
24 order.data$bin <- NA
25 order.data[order.data[,order.value] >= bin[1], 'bin'] <- 'bin1'
26 order.data[order.data[,order.value] >= bin[2] & order.data[,order.value] < bin
    [1], 'bin'] <- 'bin2'
27 order.data[order.data[,order.value] >= bin[3] & order.data[,order.value] < bin
    [2], 'bin'] <- 'bin3'
28 order.data[order.data[,order.value] >= bin[4] & order.data[,order.value] < bin
    [3], 'bin'] <- 'bin4'
29 order.data[order.data[,order.value] >= bin[5] & order.data[,order.value] < bin
    [4], 'bin'] <- 'bin5'
30
31 #converting compounds to character class and also getting compounds in each bin
32 bin.cmpd <- list()
33 for(i in c('bin1','bin2','bin3','bin4','bin5')) {
34   temp <- order.data[order.data$bin == i,]
35   bin.cmpd[[i]] <- as.character(unique(temp[,variable]))
36 }
37
38 #Getting only compounds found in one bin and not any other bins before it
39 common.cmpd <- c(bin.cmpd[[1]])
40 uni.list <- list(bin.cmpd[[1]])
41
42 for(i in 2:length(bin.cmpd)) {
43   uni.list[[i]] <- bin.cmpd[[i]][!bin.cmpd[[i]] %in% common.cmpd]
44   common.cmpd <- c(common.cmpd, uni.list[[i]])
45 }
46
47 #within bins, ordering compounds by order.value and setting variable order
48 ordered.data <- c()
49 ordered.list <- list()
50 for(i in 1:length(uni.list)) {
51   cmpd <- uni.list[[i]]
52   temp.data <- order.data[order.data[,order.var] %in% cmpd,]
53   colnames(temp.data)[colnames(temp.data)==order.value] <- 'order.val'
54   max.data <- ddply(temp.data, order.var, summarize,
55                     max.value=max(order.val))
56   temp.order <- max.data[order(max.data$max.value, decreasing=TRUE),]
57   temp.order[,order.var] <- factor(temp.order[,order.var],
58                                   levels=rev(temp.order[,order.var]),
59                                   ordered=TRUE)
60
61   temp <- data[data[,variable] %in% cmpd,]
62   temp[,variable] <- factor(temp[,variable], levels=levels(temp.order[,order.var]
63   ))
64   ordered.data <- rbind(ordered.data, temp)
65
66   temp.list <- unique(temp[,variable])
67   temp.list <- temp.list[order(temp.list, decreasing=TRUE)]
68   ordered.list[[i]] <- as.character(temp.list)
69 }
70 return(list('uni.list'=ordered.list, 'ordered.data'=ordered.data))
71
72 }

```

function-bin function.R

```

1 #####
2 #Function to bin data
3
4 #input: data = dataframe that contains the value to be represented in plot; long
    form
5 #   bin = vector of 5 values are the bin limits, in decreasing order
6 #   variable = name of column that holds the variables to be binned
7 #   value = name of column that holds values that determines the bins
8

```

```

9 #output: uni.list = class list; list of compounds in each bin
10 #
11 bin_function <- function(data, bin, variable='variable', value='value') {
12
13   #Assigning data to bins
14   data$bin <- NA
15   data[data[,value] >= bin[1], 'bin'] <- 'bin1'
16   data[data[,value] >= bin[2] & data[,value] < bin[1], 'bin'] <- 'bin2'
17   data[data[,value] >= bin[3] & data[,value] < bin[2], 'bin'] <- 'bin3'
18   data[data[,value] >= bin[4] & data[,value] < bin[3], 'bin'] <- 'bin4'
19   data[data[,value] >= bin[5] & data[,value] < bin[4], 'bin'] <- 'bin5'
20
21   #converting compounds to character class and also getting compounds in each bin
22   bin.cmpd <- list()
23   for(i in c('bin1','bin2','bin3','bin4','bin5')) {
24     temp <- data[data$bin == i,]
25     bin.cmpd[[i]] <- as.character(unique(temp[,variable]))
26   }
27
28   #Getting only compounds found in one bin and not any other bins before it
29   common.cmpd <- c(bin.cmpd[[1]])
30   uni.list <- list(bin.cmpd[[1]])
31
32   for(i in 2:length(bin.cmpd)) {
33     uni.list[[i]] <- bin.cmpd[[i]][!bin.cmpd[[i]] %in% common.cmpd]
34     common.cmpd <- c(common.cmpd, uni.list[[i]])
35   }
36
37   return(uni.list)
38 }

```

defined binning.R

```

1 #####
2 #Modified for Liquid Gold Manuscript A
3 #written for Manuscript A — function to define bins, and return which bin
4 #each compound is in.
5
6 defined_binning <- function() {
7
8   #Setting working directory
9   setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Analysis')
10
11   #loading required packages
12   require(RSQLite)
13   require(ggplot2)
14   require(plyr)
15   require(reshape2)
16   require(gridExtra)
17
18   #sourcing required functions
19   source('../Functions/function-data cleanup.R')
20   source('../Functions/function-descriptive statistics.R')
21
22   source('../Functions/function-heat map bin.R')
23
24
25   #Input data
26
27   # Choose driver
28   drv <- dbDriver('SQLite')
29
30   # Create connection (for SQLite, this is a file)
31   con <- dbConnect(drv, '../Data source files/data.db')
32
33   input <- dbGetQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
                               q.hour, q.phase, q.dose, q.treatment, q.gen_treat,

```



```

34         q.expt, d.variable, d.value
35     FROM Data d INNER JOIN Qualifier q
36     ON d.file = q.file
37     INNER JOIN Rank r
38     ON d.variable = r.Compound AND q.expt = r.expt
39     WHERE r.Confidence = 1 AND q.gen_treat = 'None'
40     AND d.variable NOT IN ('Xanthine', 'Glycerol')
41     AND q.state = 'ss'
42     AND (q.run != 8 OR q.run IS NULL)
43     AND q.donor != 'RP1'
44     AND q.expt != 'RP2'
45     ORDER BY q.sampleID;")
46
47     cmpd.info <- dbGetQuery(con, "SELECT c.Compound, c.Type, p.Pathway, m.main
48                                FROM Cmpd_type c INNER JOIN Pathway p
49                                ON c.Compound = p.Compound
50                                INNER JOIN Main_pathway m
51                                ON p.Pathway = m.path
52                                ORDER BY c.Compound")
53
54     dbDisconnect(con)
55
56     #Setting working directory
57     setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/LG Manuscript
58           A/Analysis')
59
60     #data cleanup
61
62     #Performing data cleanup to remove zero value compounds
63     #also casts data into wide format
64     clean <- data_cleanup(input, form.y=c('sampleID', 'donor', 'run', 'vessel',
65                                           'DPI', 'hour', 'phase', 'dose', 'treatment',
66                                           'gen_treat', 'expt'),
67                           form.x='variable', convert.NA=TRUE)
68
69     #make sure all compounds have a value for each donor
70     data.cast <- dcast(clean, formula=sampleID+donor+run+vessel+DPI+hour+phase+dose+
71                       treatment+gen_treat+expt~variable, value.var='value')
72
73     data.cast[is.na(data.cast)] <- 0
74
75     #melting back to long format
76     data.melt <- melt(data.cast, id.vars=c('sampleID', 'donor', 'run', 'vessel',
77                                           'DPI', 'hour', 'phase', 'dose', 'treatment',
78                                           'gen_treat', 'expt'))
79
80     #descriptive stats
81
82     #Getting descriptive stats for each donor
83     stat <- ddply(data.melt, 'donor', descriptive.stat, data.type='long', variable=
84                   variable,
85                   qualifier=qualifier, scale.data=FALSE)
86
87     #Assigning data to bins
88     bin.data <- bin_heatmap(stat, bin=c(20,10,1,0.1,0), 'variable',
89                             order.data=stat, order.var='variable', order.value='avg')
90
91     uni.list <- bin.data[['uni.list']]
92
93     #getting dataframe of variables and their bins
94     bins <- c()
95     for(i in 1:length(uni.list)) {
96         d <- data.frame(bin=i, variable=uni.list[[i]])
97         bins <- rbind(bins, d)
98     }

```

```

96 }
97
98 return(bins)
99 }

```

function-weighted compound class.R

```

1 #####
2 /
3 #Function to illustrate profile diversity
4 #Getting proportion of compound types weighted to their concentrations
5 #to be used in ddply, grouped by donor
6 custom_perc <- function(d, compd.class='Type') {
7
8   #First, getting sum of all compound concentrations for each sample
9   output <- ddply(d, 'sampleID', mutate, sample.total=sum(value, na.rm=TRUE))
10
11   #Getting mean profile totals
12   output$sample.avg <- mean(output$sample.total)
13
14   #Sum of compound classes within samples
15   output <- ddply(output, c('sampleID',compd.class), mutate,
16     sample.class=sum(value, na.rm=TRUE))
17
18   #Mean compound class
19   output <- ddply(output, compd.class, mutate, class.avg=mean(sample.class))
20
21   #Mean compound class as percent of profile total
22   output$perc <- output$class.avg/output$sample.avg*100
23
24   out <- unique(output[,c('donor', 'sample.avg',compd.class, 'class.avg', 'perc')])
25   return(out)
26 }

```

function-defined binning.R

```

1 #####
2 #####
3 #Function to get order of control data - modified for thesis
4
5 #input: data = dataframe of non-treated samples;
6 #           holds the values by which the compounds should be ordered
7
8 #output: order.list = dataframe of data to be plotted, and compounds are ordered
9 #           by a decreasing value
10 #-----
11
12 defined_binning <- function() {
13
14   wd <- getwd()
15
16   #Setting working directory
17   setwd('C:/Users/Sandi/Dropbox/Projects/Liquid Gold/R Scripts/V 3.0/Analysis')
18
19   #loading required packages
20   require(RSQLite)
21   require(ggplot2)
22   require(plyr)
23   require(reshape2)
24   require(gridExtra)
25
26   #sourcing required functions
27   source(' ../Functions/function-data cleanup.R')
28   source(' ../Functions/function-descriptive statistics.R')
29   source(' ../Functions/function-heat map bin.R')

```

```

30
31
32 #Input data


---


33 # Create connection and selecting database
34 con <- odbcConnect('MySQLDSN')
35
36 input <- sqlQuery(con, "SELECT q.sampleID, q.donor, q.run, q.vessel, q.DPI,
37       q.hour, q.phase, q.dose, q.treatment, q.gentreat, q.expt,
38       d.variable, d.value, t.tax_class
39 FROM liquidgold.data d INNER JOIN liquidgold.qualifier q
40 ON d.file = q.file
41 INNER JOIN liquidgold.rank r
42 ON d.variable = r.Compound AND q.expt = r.expt
43 INNER JOIN general.tax t
44 ON d.variable=t.syn
45 WHERE r.Confidence = 1 AND q.gentreat = 'None'
46 AND (q.run != 8 OR q.run IS NULL)
47 AND d.variable NOT IN ('Xanthine', 'Glycerol')
48 AND q.state = 'ss'
49 AND q.donor != 'RP1'
50 AND q.expt != 'RP2'
51 ORDER BY q.sampleID;")
52
53 odbcClose(con)
54
55 #data cleanup


---


56
57 #Performing data cleanup to remove zero value compounds
58 #also casts data into wide format
59 clean <- data_cleanup(input, convert.NA=TRUE)
60
61 #make sure all compounds have a value for each donor
62 data.cast <- dcast(clean, formula=sampleID+donor+run+vessel+DPI+hour+phase+dose+
63       treatment+gentreat+expt~variable, value.var='value')
64
65 data.cast[,12:ncol(data.cast)][is.na(data.cast[,12:ncol(data.cast)])] <- 0
66
67 #melting back to long format
68 data.melt <- melt(data.cast, id.vars=c('sampleID', 'donor', 'run', 'vessel',
69       'DPI', 'hour', 'phase', 'dose', 'treatment',
70       'gentreat', 'expt'))
71
72 data.melt <- unique(merge(data.melt, clean[,c('tax_class', 'variable')],
73       by='variable'))
74 #separating out qualifier and data
75 qualifier <- data.melt[,1:11]
76 x.data <- data.melt[,12:13]
77
78 #descriptive stats


---


79 #Getting descriptive stats for each donor
80 stat <- ddply(data.melt, 'donor', descriptive.stat, data.type='long',
81       variable='variable')
82
83 #Binning


---


84 #Assigning data to bins
85 bin.data <- bin.heatmap(stat, bin=c(20,10,1,0.1,0), 'variable',
86       order.data=stat, order.var='variable', order.value='avg')
87
88 uni.list <- bin.data[['uni.list']]
89 ordered.data <- bin.data[['ordered.data']]

```

```
90 |  
91 | var <- levels(ordered.data$variable)  
92 | var <- factor(var, levels=var, ordered=TRUE)  
93 |  
94 | for(i in 1:length(uni.list)){  
95 |   uni.list[[i]] <- factor(uni.list[[i]], levels=var)  
96 | }  
97 |  
98 | #setting working directory back  
99 | setwd <- wd  
100 | return(uni.list)  
101 | }
```